

Tail Analysis without Parametric Models: A Worst-case Perspective

Henry Lam

Department of Industrial and Operations Engineering, University of Michigan, Ann Arbor, MI 48109, khlam@umich.edu

Clementine Mottet

Department of Mathematics and Statistics, Boston University, Boston, MA 02215, cmottet@bu.edu

A common bottleneck in evaluating extremal performance measures is that, due to their very nature, tail data are often very limited. The conventional approach selects the best probability distribution from tail data using parametric fitting, but the validity of the parametric choice can be difficult to verify. This paper describes an alternative based on the computation of worst-case bounds under the geometric premise of tail convexity, a feature shared by all common parametric tail distributions. We characterize the optimality structure of the resulting optimization problem, and demonstrate that the worst-case convex tail behavior is in a sense either extremely light-tailed or extremely heavy-tailed. We develop low-dimensional nonlinear programs that distinguish between the two cases and compute the worst-case bound. We numerically illustrate how the proposed approach can give more reliable performances than conventional parametric methods.

Key words: tail modeling, robust analysis, nonparametric

1. Introduction

Modeling extreme behaviors is a fundamental task in analyzing and managing risk. As the earliest applications, hydrologists and climatologists study historical data of sea levels and air pollutants to estimate the risk of flooding and pollution (Gumbel (2012)). In non-life or casualty insurance, insurers rely on accurate prediction of large losses to price and manage insurance policies (McNeil (1997), Beirlant and Teugels (1992), Embrechts et al. (1997)). Relatedly, financial managers estimate risk measures of portfolios to safeguard losses (Glasserman and Li (2005), Glasserman et al. (2007, 2008)). In engineering, measurement of system reliability often involves modeling the tail behaviors of individual components' failure times (Nicola et al. (1993), Heidelberger (1995)).

Despite its importance in various disciplines, tail modeling is an intrinsically difficult task because, by their own nature, tail data are often very limited. Consider these two examples:

EXAMPLE 1 (ADOPTED FROM MCNEIL (1997)). *There were 2,156 Danish fire losses over one million Danish Krone (DKK) from 1980 to 1990. The empirical cumulative distribution function (ECDF) and the histogram (in log scale) are plotted in Figure 1. For a concrete use of the data, an insurance company might be interested in pricing a high-excess contract with reinsurance, which has a payoff of $X - 50$ (in million DKK) when $50 < X \leq 200$, 150 when $X > 200$, and 0 when $X \leq 50$, where X is the loss amount (the marks 50 and 200 are labeled with vertical lines in Figure 1). Pricing this contract would require, among other information, $E[\text{payoff}]$. However, only seven data points are above 50 (the loss amount above which the payoff is non-zero).*

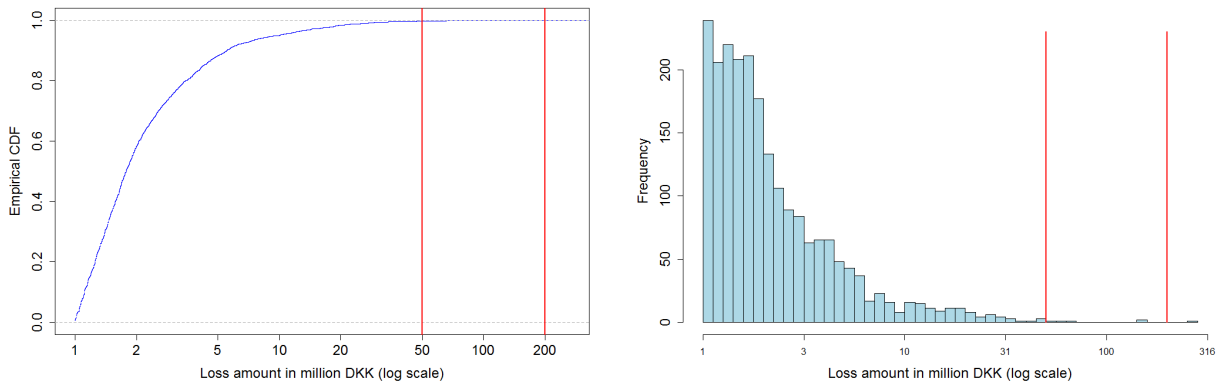


Figure 1: ECDF and histogram for Danish fire losses from 1980 to 1990

EXAMPLE 2. *A more extreme situation is a synthetic data set of size 200 generated from an unknown distribution, whose histogram is shown in Figure 2. Suppose the quantity of interest is $P(4 < X < 5)$. This appears to be an ill-posed problem since the interval $[4, 5]$ has no data at all. This situation is not uncommon when in any application one tries to extrapolate the tail with a small sample size.*

The purpose of this paper is to develop a theoretically justified methodology to estimate tail-related quantities of interest such as those depicted in the examples above. This requires drawing information properly from data not in the tail. We will illustrate how to do this and revisit the two examples later with numerical performance of our method.

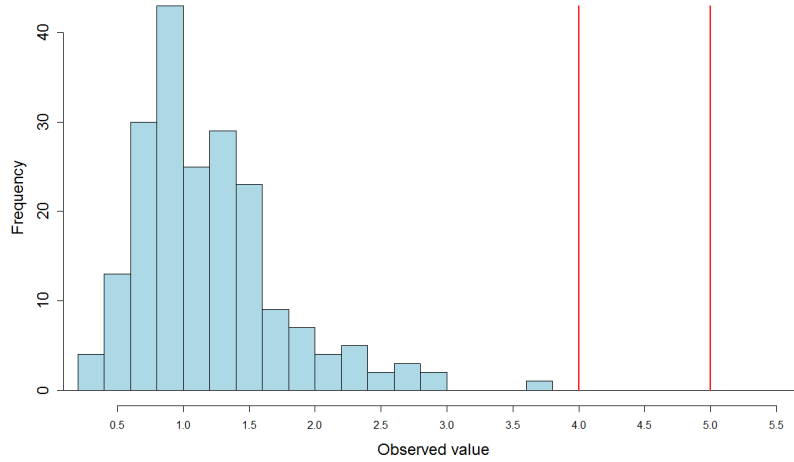


Figure 2 Histogram of a synthetic data set with sample size 200

2. Our Approach and Main Contributions

We adopt a nonparametric approach. Rather than fitting a tail parametric curve when there can be few or zero observations in the tail region, we base our analysis on the geometric premise that the tail density is convex. We emphasize that this condition is satisfied by *all* known parametric distributions (e.g. normal, lognormal, exponential, gamma, Weibull, Pareto etc.). For this reason we believe it is a natural and minimal assumption to make.

In any given problem, there can be potentially infinitely many feasible candidates of convex tails. The central idea of our method is a worst-case characterization. Formally, given information on the non-tail part of the distribution and a target quantity of interest (e.g., $P(4 < X < 5)$ in Example 2), we aim to find a convex tail, consistent with the non-tail part, that gives rise to the worst-case value of the target (e.g., the largest possible value of $P(4 < X < 5)$). This value serves as a tight bound for the target that is robust with respect to the ambiguity of the tail, without using any particular tail knowledge other than our a priori assumption of convexity.

Our proposed approach requires solving an optimization over a potentially infinite-dimensional space of convex tails. As our key contributions, we show that this problem has a very simple optimality structure, and find its solution via low-dimensional nonlinear programs. In particular:

1. We characterize the worst-case tail behavior under the tail convexity condition. We show that the worst-case tail, for *any* bounded target quantity of interest, is in a sense either *extremely light-tailed* or *extremely heavy-tailed*. Both cases can be characterized by piecewise linear densities, the distinction being whether the pieces form a bounded support distribution or lead to probability masses that escape to infinity.

2. We provide efficient algorithms to distinguish between the two cases above, and to solve for the optimal distribution in each case. For a large class of objectives, the algorithm requires at most a two-dimensional nonlinear program.

Our approach outputs statistically valid worst-case bounds when integrating with confidence estimates drawn from the non-tail portion of the data. This approach uses the convexity assumption to get around the difficulty faced by conventional parametric methods (discussed in detail in the next section) in directly estimating the tail curve, by effectively mitigating the estimation burden to the central part of the density curve where more data are available. However, we pay the price of conservativeness: our method can generate a worst-case bound that is over-pessimistic. We therefore believe it is most suitable for small sample size, when a price of conservativeness is unavoidable in trading with statistical validity.

The remainder of this paper is organized as follows. Section 3 discusses some previous techniques and reviews the relevant literature. Section 4 presents our formulation and results for an abstract setting. Section 5 studies the numerical solution algorithm. Section 6 focuses on integrating these results with data. Section 7 shows some numerical illustration. Section 8 concludes and discusses future work. Some auxiliary theorems and proofs are left to the Appendix.

3. Related Work

3.1. Overview of Common Tail-fitting Techniques

As far as we know, all existing techniques for modeling extreme events are parametric-based, in the sense that a “best” parametric curve is chosen and the parameters are fit to the tail data. The classic text of Hogg and Klugman (2009) provides a comprehensive discussion on the common

choices of parametric tail densities. While exploratory data analysis, such as quantile plots and mean excess plots, can provide guidance regarding the class of parametric curves to use (such as heavy, middle or light tail), this approach is limited by its reliance on a large amount of data in the tail and subjectivity in the choice of parametric curve.

Beyond the goodness-of-fit approach, there are two widely used results on the parametric choice that is provably suitable for extreme values. The Fisher-Tippett-Gnedenko Theorem (Fisher and Tippett (1928), Gnedenko (1943)) postulates that the sample maxima, after suitable scaling, must converge to a generalized extreme value (GEV) distribution, given that it converges at all to some non-degenerate distribution. This result is useful if the data are known to derive from the maximum of some distributions. For instance, environmental data on sea level and river heights are often collected as annual maxima (Davison and Smith (1990)), and in this scenario it is sensible to fit the GEV distribution. In other scenarios, the data have to be pre-divided into blocks and blockwise maxima have to be taken in order to apply GEV, but this blockwise approach is statistically wasteful (Embrechts et al. (2005)).

The Pickands-Balkema-de Haan Theorem (Pickands III (1975), Balkema and De Haan (1974)) does not require data to come from maxima. Rather, the theorem states that the excess losses over thresholds converge to a generalized Pareto distribution (GPD) as the thresholds approach infinity, under the same conditions as the Fisher-Tippett-Gnedenko Theorem. The Pickands-Balkema-de Haan theorem provides a solid mathematical justification for using GPD to fit the tail portion of data (McNeil (1997), Embrechts et al. (2005)). Fitting GPD can be done by well-studied procedures such as maximum likelihood estimation (Smith (1985)), and the method of probability-weighted moments (Hosking and Wallis (1987)). The Hill estimator (Hill et al. (1975), Davis and Resnick (1984)) is also a widely used alternative.

Despite the attraction and frequent usage, fitting GPD suffers from two pitfalls: First, there is no convergence rate result that tells how high a threshold should be for the GPD approximation to be valid (e.g. McNeil (1997)). Hence, picking the threshold is an ad hoc task in practice. Second,

and more importantly, even if the threshold chosen is sufficiently high for the approximation to hold, a large amount of data above it is needed to accurately estimate the parameters in GPD. In our two examples, especially Example 2, this is plainly impossible.

3.2. Related Literature on our Methodology

Our mathematical formulation and techniques are related to two lines of literature. The use of convexity and other shape constraints (such as log-concavity) have appeared in density estimation (Cule et al. (2010), Seregin and Wellner (2010), Koenker and Mizera (2010)) and convex regression (Seijo et al. (2011), Hannah and Dunson (2013), Lim and Glynn (2012)) in statistics. A major reason for using convexity in these statistical problems is the removal of tuning parameters, such as bandwidth, as required by other methods such as the use of kernel.

The second line of related literature is optimization over probability distributions, which have appeared in decision analysis (Smith (1995), Bertsimas and Popescu (2005), Popescu (2005)), robust control theory (Iyengar (2005), El Ghaoui and Nilim (2005), Petersen et al. (2000), Hansen and Sargent (2008)), distributionally robust optimization (Delage and Ye (2010), Goh and Sim (2010)), and stochastic programming (Birge and Wets (1987), Birge and Dula (1991)). The typical formulation involves optimization of some objective governed by a probability distribution that is partially specified via constraints like moments (Karr (1983), Winkler (1988)) and statistical distances (Ben-Tal et al. (2013)). Our formulation differs from these studies by its pertinence to tail modeling (i.e., knowledge of certain regions of the density, but none beyond it). Among all the previous works, only Popescu (2005) has considered convex density assumption, as an instance of a proposed class of geometric conditions that are added to moment constraints. While the result bears similarity to ours in that a piecewise linearity structure shows up in the solution, our qualitative classification of the tail, the solution techniques, and the formulation in integrating with data all differ from the semidefinite programming approach in Popescu (2005).

4. Abstract Formulation and Results

We begin by considering an abstract formulation assuming full information on the distribution up to some threshold, and no information beyond. The next sub-sections give the details.

4.1. Formulation

Consider a continuous probability distribution on \mathbb{R} whose density exists and is denoted by $f(x)$. We assume that f is known up to a certain large threshold, say $a \in \mathbb{R}$. The goal is to extrapolate f .

We impose the assumption that $f(x)$, for $x \geq a$, is convex. Figure 3 shows an example of an $f(x)$ known up to a , and Figures 4 and 5 each show an example of convex and non-convex extrapolation. Observe that the convex tail assumption excludes any “surprising” bumps (and falls) in the density curve.

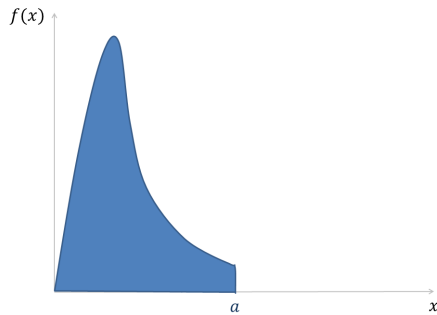


Figure 3: A probability density $f(x)$ known up to a threshold a

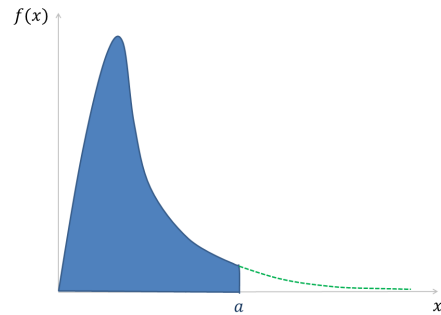


Figure 4: An example of convex tail extrapolation

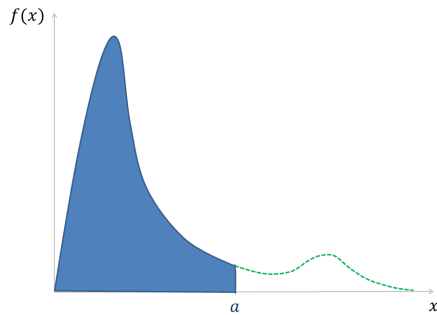


Figure 5: An example of non-convex tail extrapolation

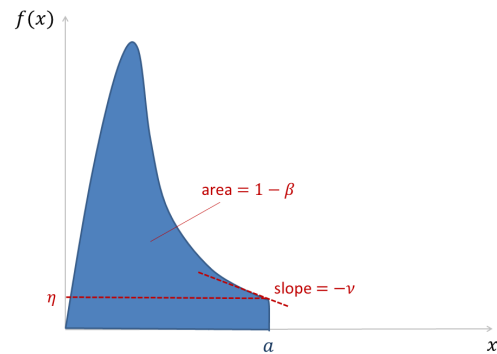


Figure 6: The parameters η, ν, β

Now suppose we are given a target objective or performance measure $E[h(X)]$, where $E[\cdot]$ denotes the expectation under f , and $h: \mathbb{R} \rightarrow \mathbb{R}$ is a bounded function in X . The goal is to calculate the worst-case value of $E[h(X)]$ under the assumption that f is convex beyond a . That is, we want

to obtain $\max E[h(X)] = \int_{-\infty}^{\infty} h(x)f(x)dx$ where the maximization is over all convex $f(x), x \geq a$ such that it satisfies the properties of a probability density function. We assume that the density is left-differentiable at a , so that a convex extrapolation at a can be suitably defined. For the formulation, we need three constants extracted from $f(x), x < a$, which we denote as $\eta, \nu, \beta > 0$ respectively:

1. η is the value of the density f at a , i.e. $f(a) = \eta$.
2. $-\nu$ is the left derivative of f at a , i.e. $f'_-(a) = -\nu$. We impose the condition that the right derivative $f'_+(a) \geq f'_-(a) = -\nu$. Note that, since f is convex (and bounded) on $[a, \infty)$, its one-sided derivative exists everywhere on $[a, \infty)$.
3. β is the tail probability at a . Since f is known up to a , $\int_{-\infty}^a f(x)$ is known to be equal to some number $1 - \beta$, and $\int_a^{\infty} f(x)dx$ must equal β .

Figure 6 illustrates these quantities. For $\eta, \nu, \beta > 0$, our formulation can be written as

$$\begin{aligned} \max_f \quad & \int_a^{\infty} h(x)f(x)dx \\ \text{subject to} \quad & \int_a^{\infty} f(x)dx = \beta \end{aligned} \tag{1a}$$

$$f(a) = f(a+) = \eta \tag{1b}$$

$$f'_+(a) \geq -\nu \tag{1c}$$

$$f \text{ convex for } x \geq a \tag{1d}$$

$$f(x) \geq 0 \text{ for } x \geq a \tag{1e}$$

Note that we have set our objective to be $E[h(X); X \geq a]$, since $E[h(X); X < a]$ is completely known in this setting. Here $f(a+)$ denotes the right-limit at a , and $f(a) = f(a+)$ means that f is right-continuous at a , implying a continuous extrapolation at a .

4.2. Optimality Characterization

The solution structure of (1) turns out to be extremely simple and is characterized by either one of two closely related cases (focusing on the region $x \geq a$). Let $\mathcal{C}^+[a, \infty)$ denote the class of non-negative continuous functions on $[a, \infty)$. Let

$$\mathcal{PL}_m^+[a, \infty) = \{f \in \mathcal{C}^+[a, \infty) : f(x) = c_j + d_j x \text{ for } x \in [y_{j-1}, y_j], j = 1, \dots, m,$$

where $a = y_0 \leq y_1 \leq \dots \leq y_m < \infty$, $c_j, d_j \in \mathbb{R}$, and $f(x) = 0$ for $x > y_m$

be the set of all non-negative, continuous and piecewise linear functions on $[a, \infty)$ that have at most m line segments before vanishing. We have:

THEOREM 1. *Suppose h is measurable and bounded. Consider optimization (1). If it is feasible, then either*

1. *An optimal solution f^* exists, where $f^* \in \mathcal{PL}_3^+[a, \infty)$.*
2. *An optimal solution does not exist. There exists a sequence $\{f^{(k)} \in \mathcal{PL}_3^+[a, \infty) : k \geq 1\}$, each $f^{(k)}$ feasible for (1), such that $\int_a^\infty h(x)f^{(k)}(x)dx \rightarrow Z^*$ as $k \rightarrow \infty$, where Z^* is the optimal value of (1). Moreover, let $\{c_3^{(k)} + d_3^{(k)}x : x \in [y_2^{(k)}, y_3^{(k)}]\}$ be the last line segment of $f^{(k)}$. We have $y_3^{(k)} \nearrow \infty$ and $d_3^{(k)} \searrow 0$ as $k \rightarrow \infty$.*

The proof of Theorem 1 is discussed in the next sub-sections. Note that f^* in the first case in Theorem 1 is a continuous piecewise linear density, and consequently has bounded support. In the second case, as $k \rightarrow \infty$, the sequence $\{f^{(k)} : k \geq 1\}$ has unboundedly increasing support endpoint ($y_3^{(k)} \nearrow \infty$), and its last line segment gets closer and more parallel to the horizontal axis ($d_3^{(k)} \searrow 0$). This sequence possesses a pointwise limit, but the limit is not a valid density and has a probability mass that “escapes” to positive infinity.

Figures 7 and 8 show the tail behaviors for the two cases above. A bounded support density in the first case possesses the lightest possible tail behavior. The second case, on the other hand, can be interpreted as an extreme heavy-tail. Compare the sequence $f^{(k)}$ with a given arbitrary density. Given any fixed large enough x on the real line, as k grows, the decay rate of $f^{(k)}$ at the point x is eventually slower than that of the given density. Since a slower decay rate is the characteristic of a fatter tail, the behavior implied by $f^{(k)}$ in a sense captures the heaviest possible tail.

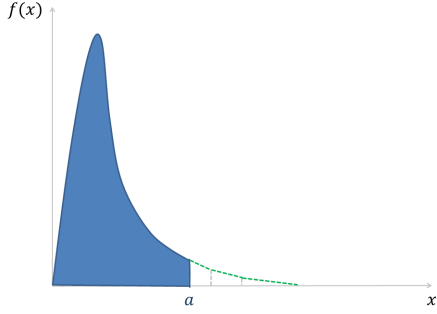


Figure 7: Behavior of an optimal light-tailed extrapolation

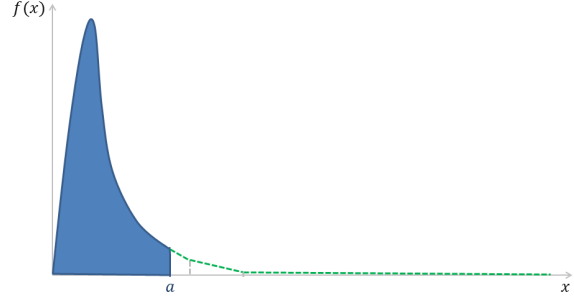


Figure 8: Behavior of an element in an optimal heavy-tailed extrapolation sequence

4.3. Main Mathematical Developments

This section presents the mathematical argument for Theorem 1. This development will also help construct a solution algorithm in Section 5. We divide the argument into two parts. First we establish an equivalence of (1) to a moment-constrained optimization problem under a different probability space. Second, we characterize the solution of this moment-constrained problem, which can then be converted to the solution of (1)

We define some notations. Let \mathbb{R}^+ and \mathbb{R}^- be the non-negative and non-positive real axis. Denote $\mathcal{P}(\mathcal{M})$ as the set of all probability measures on a measurable space \mathcal{M} equipped with the Borel σ -field. Let $\mathcal{S}_l = \{(p_1, \dots, p_l) \in (\mathbb{R}^+)^l : \sum_{i=1}^l p_i = 1\}$ be the l -dimensional probability simplex. Let $\delta(\cdot)$ be the Dirac measure. Denote $\mathcal{P}_n(\mathcal{M})$ as the set of all finite support distributions on \mathcal{M} with at most n support points, i.e. each $\mathbb{P} \in \mathcal{P}_n(\mathcal{M})$ has masses $p_1, p_2, \dots, p_l \in \mathcal{S}_l$ on points $x_1, \dots, x_l \in \mathcal{M}$, where $1 \leq l \leq n$, defined such that $\mathbb{P} = \sum_{i=1}^n p_i \delta(x_i)$. For simplicity, since any $\mathbb{P} \in \mathcal{P}_n(\mathcal{M})$ can be represented by the support points $(x_1, \dots, x_n) \in \mathcal{M}^n$ (some possibly identical) and $(p_1, \dots, p_n) \in \mathcal{S}_n$, we sometimes write $\mathbb{P} \sim (x_1, \dots, x_n, p_1, \dots, p_n)$ for a given $\mathbb{P} \in \mathcal{P}_n(\mathcal{M})$. Moreover, we use the notation $\mathbb{E}[\cdot]$ to denote the associated expectation under \mathbb{P} .

For convenience, denote $\mathcal{P}^+ = \mathcal{P}(\mathbb{R}^+)$ as the set of all probability measures concentrated on \mathbb{R}^+ , and $\mathcal{P}_n^+ = \mathcal{P}_n(\mathbb{R}^+)$ the corresponding set of measures with at most n support points. The measurability of h is assumed throughout the rest of the exposition.

4.3.1. Equivalence to Moment-constrained Optimization

We first reformulate (1) as:

LEMMA 1. *Formulation (1) is equivalent to*

$$\begin{aligned} \max_f \quad & \int_a^\infty h(x)f(x)dx \\ \text{subject to} \quad & \int_a^\infty f(x)dx = \beta \end{aligned} \tag{2a}$$

$$f(a) = \eta \tag{2b}$$

$$f'_+(x) \text{ exists and is non-decreasing and right-continuous for } x \geq a \tag{2c}$$

$$-\nu \leq f'_+(x) \leq 0 \text{ for } x \geq a \tag{2d}$$

$$f'_+(x) \rightarrow 0 \text{ as } x \rightarrow \infty \tag{2e}$$

$$f(x) = \int_a^x f'_+(t)dt + \eta \text{ for } x \geq a \tag{2f}$$

Proof of Lemma 1. The proof uses several elementary results from convex analysis. See Appendix EC.1 for details. \square

As a key step, we show the equivalence of (2) to a moment-constrained program, by identifying the decision variable as $f'_+(x)$ via a one-to-one map with a probability distribution function. Let

$$H(x) = \int_0^x \int_0^u h(v+a)dvdu \tag{3}$$

and

$$\mu = \frac{\eta}{\nu} \text{ and } \sigma = \frac{2\beta}{\nu} \tag{4}$$

where $\mu, \sigma > 0$ since we have assumed $\beta, \eta, \nu > 0$. Our result is:

THEOREM 2. *Suppose h is bounded. The optimal value of (2) is equal to that of*

$$\begin{aligned} \max_{\mathbb{P}} \quad & \nu \mathbb{E}[H(X)] \\ \text{subject to} \quad & \mathbb{E}[X] = \mu \\ & \mathbb{E}[X^2] = \sigma \\ & \mathbb{P} \in \mathcal{P}^+ \end{aligned} \tag{5}$$

Here the decision variable is a probability measure $\mathbb{P} \in \mathcal{P}^+$, and $\mathbb{E}[\cdot]$ is the corresponding expectation. Moreover, there is a one-to-one correspondence between the feasible solutions to (2) and (5), given by $f'_+(x+a) = \nu(p(x) - 1)$ for $x \in \mathbb{R}^+$, where f'_+ is the right derivative of a feasible solution f of (2) such that $f(x) = \int_a^x f'_+(t)dt + \eta$ for $x \geq a$, and p is a probability distribution function that is associated with a feasible probability measure over \mathbb{R}^+ in (5).

Proof of Theorem 2. The key of the proof uses integration by parts and an explicit construction of a linear transformation between f'_+ and a probability distribution function p . See Appendix EC.1 for details. \square

4.3.2. Further Reduction and Optimality Characterization Next we characterize the optimality structure for (5), a generalized moment problem in the form of an infinite-dimensional linear program. Using existing terminology, we call an optimization program *consistent* if there exists a feasible solution, and *solvable* if there exists an optimal solution.

For convenience, denote $OPT(\mathcal{D})$ as the program

$$\begin{aligned} \max_{\mathbb{P}} \quad & \nu \mathbb{E}[H(X)] \\ \text{subject to} \quad & \mathbb{E}[X] = \mu \\ & \mathbb{E}[X^2] = \sigma \\ & \mathbb{P} \in \mathcal{D} \end{aligned}$$

where H, μ, σ are defined in (3) and (4), and \mathcal{D} is a collection of probability measures on \mathbb{R} . For example, program (5) is denoted as $OPT(\mathcal{P}^+)$. Moreover, let $Z(\mathbb{P}) = \nu \mathbb{E}[H(X)]$ be the objective function of $OPT(\mathcal{D})$ in terms of \mathbb{P} . We have:

THEOREM 3. *Program (5), or equivalently $OPT(\mathcal{P}^+)$, has the same optimal value as $OPT(\mathcal{P}_3^+)$.*

Proof of Theorem 3. Follows from a classical result on the extreme points of moment sets. See Appendix EC.1. \square

Next we derive some properties regarding the optimality of $OPT(\mathcal{P}_3^+)$:

PROPOSITION 1. Consider $OPT(P_3^+)$ that is consistent. The optimal value Z^* is either achieved at some $\mathbb{P}^* \in \mathcal{P}_3^+$, or there exists a sequence of feasible $\mathbb{P}^{(k)} \in \mathcal{P}_3^+$ such that $Z(\mathbb{P}^{(k)}) \rightarrow Z^*$. In the second case, each $\mathbb{P}^{(k)} \sim (x_1^{(k)}, x_2^{(k)}, x_3^{(k)}, p_1^{(k)}, p_2^{(k)}, p_3^{(k)})$, such that either $(x_1^{(k)}, x_2^{(k)}, x_3^{(k)}, p_1^{(k)}, p_2^{(k)}, p_3^{(k)}) \rightarrow (x_1^*, x_2^*, \infty, p_1^*, p_2^*, 0)$ for some $x_1^*, x_2^* \in \mathbb{R}^+$ (possibly identical) and $(p_1^*, p_2^*) \in \mathcal{S}_2$, or $(x_1^{(k)}, x_2^{(k)}, x_3^{(k)}, p_1^{(k)}, p_2^{(k)}, p_3^{(k)}) \rightarrow (x_1^*, \infty, \infty, 1, 0, 0)$ for some $x_1^* \in \mathbb{R}^+$.

Proof of Proposition 1. See Appendix EC.1. □

We are now ready to show Theorem 1:

Proof of Theorem 1. Convert the original optimization (1) into (5) by Lemma 1 and Theorem 2.

If (5) is consistent, then, by Theorem 3, its optimal value is attained by the two cases in Proposition 1. Note that that any solution $\mathbb{P} \in \mathcal{P}_3[0, \infty)$ represented by $(x_1, x_2, x_3, p_1, p_2, p_3)$ (with potentially identical x_i 's) admits one-to-one correspondence with a solution f in (1), via $f'_+(x+a) = \nu(p(x) - 1)$ in Theorem 2, giving

$$f'_+(x) = \begin{cases} -\nu & \text{for } a \leq x < x_1 + a \\ -\nu(1 - p_1) & \text{for } x_1 + a \leq x < x_2 + a \\ -\nu(1 - p_1 - p_2) & \text{for } x_2 + a \leq x < x_3 + a \\ 0 & \text{for } x_3 + a \leq x \end{cases}$$

and hence

$$f(x) = \begin{cases} \eta - \nu(x - a) & \text{for } a \leq x \leq x_1 + a \\ \eta - \nu x_1 - \nu(1 - p_1)(x - a - x_1) & \text{for } x_1 + a \leq x \leq x_2 + a \\ \eta - \nu x_1 - \nu(1 - p_1)(x_2 - x_1) - \nu(1 - p_1 - p_2)(x - a - x_2) & \text{for } x_2 + a \leq x \leq x_3 + a \\ 0 & \text{for } x_3 + a \leq x \end{cases} \quad (6)$$

The first case in Proposition 1 thus concludes Part 1 of Theorem 1. In the second case in Proposition 1, $x_3^{(k)} \rightarrow \infty$ and $p_3^{(k)} \rightarrow 0$ so that $1 - p_1^{(k)} - p_2^{(k)} \rightarrow 0$. Using (6), we conclude Part 2 of Theorem 1.

□

We close this section with two results. First is on the consistency of programs (1) and (5):

LEMMA 2. *Program (5) is consistent if and only if $\sigma \geq \mu^2$. Correspondingly, program (1) is consistent if and only if $\eta^2 \leq 2\beta\nu$. When $\sigma = \mu^2$, (5) has only one feasible solution given by $\delta(\mu)$. Correspondingly, when $\eta^2 = 2\beta\nu$, (1) has only one feasible solution given by $f(x) = \eta - \nu(x - a)$ for $x \geq a$.*

Proof of Lemma 2. See Appendix EC.2. □

Graphically, $\eta^2 > 2\beta\nu$ implies that β is smaller than the area under the straight line starting from the point (a, η) down to the x -axis with slope $-\nu$. Hence no convex extrapolation can be drawn under this condition.

Next, we show that the boundedness assumption on h is nearly essential, in the sense that any polynomially growing h leads to an infinite optimal value for (1):

PROPOSITION 2. *Suppose $\eta^2 < 2\beta\nu$ and $h(x) = \Omega(x^\epsilon)$ as $x \rightarrow \infty$ for some $\epsilon > 0$. The optimal value of (1) is ∞ .*

Proof of Proposition 2. The proof explicitly constructs a sequence of feasible solutions that lead to exploding objective values. See Appendix EC.1. □

5. Optimization Procedure for Quasi-concave Objectives

This section develops a numerical solution algorithm for our worst-case optimization presented in Section 4. In building our algorithm, we focus on h that satisfies the following stronger assumption, which covers many natural scenarios including the two examples in the Introduction.

ASSUMPTION 1. *The function $h : \mathbb{R} \rightarrow \mathbb{R}^+$ is bounded, and is non-decreasing in $[a, c)$ and non-increasing in (c, ∞) for some constant $a \leq c \leq \infty$ (i.e. c can possibly be ∞).*

Assumption 1 implies that h is quasi-concave. The non-negativity of h is assumed without loss of generality when applied to optimization (1). Because h is bounded, one can always add a sufficiently large constant, say C , to make h non-negative. Note that we have $E[h(X); X \geq a] = E[h(X) + C; X \geq a] - CP(X \geq a) = E[h(X) + C; X \geq a] - C\beta$, and so one can solve $E[h(X) + C; X \geq a]$ and recover $E[h(X); X \geq a]$.

We impose an additional mild regularity assumption:

ASSUMPTION 2. *The limit*

$$\lambda = \lim_{x \rightarrow \infty} \frac{H(x)}{x^2} \quad (7)$$

where H is defined in (3), exists and is finite.

Note that when h is bounded, $H(x) = O(x^2)$ as $x \rightarrow \infty$, and $\limsup_{x \rightarrow \infty} H(x)/x^2 < \infty$. The essence of Assumption 2 is on the existence of the limit.

Under Assumption 2, denote

$$W(x_1) = \nu \left(\frac{\sigma - \mu^2}{\sigma - 2\mu x_1 + x_1^2} H(x_1) + \frac{(\mu - x_1)^2}{\sigma - 2\mu x_1 + x_1^2} H\left(\frac{\sigma - \mu x_1}{\mu - x_1}\right) \right) \quad (8)$$

with $W(\mu) := \nu(H(\mu) + \lambda(\sigma - \mu^2))$, where μ and σ are defined in (4). We have the following strengthened version of Theorem 1:

THEOREM 4. *Under Assumption 1,*

1. *The conclusions of Theorem 1 hold with $\mathcal{PL}_3^+[a, \infty)$ replaced by $\mathcal{PL}_2^+[a, \infty)$.*
2. *Suppose $\eta^2 < 2\beta\nu$ and Assumption 2 holds additionally. The optimal value of (1) is given by*

$$\max_{x_1 \in [0, \mu]} W(x_1).$$

3. *Suppose $\eta^2 < 2\beta\nu$ and Assumption 2 holds additionally. If there exists $x_1^* \in \operatorname{argmax}_{x_1 \in [0, \mu]} W(x_1)$ such that $x_1^* \in [0, \mu)$, then an optimal solution to (1) is given by*

$$f^*(x) = \begin{cases} \eta - \nu(x - a) & \text{for } a \leq x \leq x_1^* + a \\ \eta - \nu x_1^* - \nu \frac{(\mu - x_1^*)^2}{\sigma - 2\mu x_1^* + x_1^{*2}} (x - a - x_1^*) & \text{for } x_1^* + a \leq x \leq \frac{\sigma - \mu x_1^*}{\mu - x_1^*} + a \\ 0 & \text{for } \frac{\sigma - \mu x_1^*}{\mu - x_1^*} + a \leq x \end{cases}$$

Otherwise, there exists a sequence of feasible solutions $f^{(k)}$ with $\int_a^\infty h(x) f^{(k)}(x) dx \rightarrow Z^*$, where Z^* is the optimal value of (1). $f^{(k)} \rightarrow f^*$ pointwise where

$$f^*(x) = \begin{cases} \eta - \nu(x - a) & \text{for } a \leq x \leq \mu + a \\ 0 & \text{for } \mu + a \leq x \end{cases}$$

The second case can occur only when $\lambda > 0$.

Part 1 of Theorem 4 simplifies the search space of densities in (1) from three to two linear segments. Because of this simplification, solving (1) reduces to finding the first kink of the optimal density (or sequence of densities), equivalently the first support point of the reformulation (5). This can be done by a one-dimensional line search $\max_{x_1 \in [0, \mu]} W(x_1)$ in Part 2 of the theorem.

Part 3 of Theorem 4 describes how to distinguish between the light- and heavy-tail cases in Theorem 1 by looking at the location of x_1^* . The former case occurs when there exists a x_1^* in $[0, \mu)$, and the latter occurs otherwise. Note that $f^*(x) = 0$, $x \geq \mu + a$ in the pointwise limit of $f^{(k)}$ in Part 3 of Theorem 4 is a consequence of the last line segment of $f^{(k)}$ getting increasingly closer and more parallel to the x -axis.

Algorithm 1 summarizes the procedure for obtaining the optimal value of (1).

Algorithm 1: Procedure for finding the optimal value of (1)

Inputs:

1. The function h that satisfies Assumptions 1 and 2.
2. The parameters $\beta, \eta, \nu > 0$.

Procedure:

1. If $\eta^2 > 2\beta\nu$, there is no feasible solution.
 2. If $\eta^2 = 2\beta\nu$, the optimal value is $\nu H(\mu)$.
 3. If $\eta^2 < 2\beta\nu$, the optimal value is given by $\max_{x_1 \in [0, \mu]} W(x_1)$.
-

The rest of this section provides the developments for proving Theorem 4. First we introduce the following condition:

ASSUMPTION 3. H is convex and H' satisfies a convex-concave property, i.e. $H'(x)$ is convex for $x \in (0, c)$ and concave for $x \in (c, \infty)$, for some $0 \leq c \leq \infty$.

With Assumption 3, Theorem 3 can be strengthened to:

PROPOSITION 3. Under Assumption 3, $OPT(\mathcal{P}^+)$ has the same optimal value as $OPT(\mathcal{P}_2^+)$.

Proof of Proposition 3. See Appendix EC.2. □

This allows us to focus on one of the support points of $OPT(\mathcal{P}_2^+)$ in the solution scheme, leading to the following proposition:

PROPOSITION 4. *Under Assumptions 2 and 3, consider $OPT(\mathcal{P}_2^+)$ with $\sigma > \mu^2$ and let Z^* be its optimal value.*

1. *If there exists an optimal solution in \mathcal{P}_2^+ , then this solution has distinct support points and is represented by $(x_1^*, x_2^*, p_1^*, p_2^*)$, where $x_1^* \in \operatorname{argmax}_{x_1 \in [0, \mu]} W(x_1)$ and*

$$x_2^* = \frac{\sigma - \mu x_1^*}{\mu - x_1^*}, \quad p_1^* = \frac{\sigma - \mu^2}{\sigma - 2\mu x_1^* + x_1^{*2}}, \quad p_2^* = \frac{(\mu - x_1^*)^2}{\sigma - 2\mu x_1^* + x_1^{*2}} \quad (9)$$

Moreover, $Z^ = \max_{x_1 \in [0, \mu]} W(x_1)$.*

2. *If there does not exist an optimal solution, then there must exist a sequence $\mathbb{P}^{(k)} \sim (x_1^{(k)}, x_2^{(k)}, p_1^{(k)}, p_2^{(k)}) \rightarrow (\mu, \infty, 1, 0)$. Moreover, $Z^* = \nu(H(\mu) + \lambda(\sigma - \mu^2))$.*

3. $Z^* = \max_{x_1 \in [0, \mu]} W(x_1)$

Proof of Proposition 4. See Appendix EC.2. □

The following corollary provides a simple sufficient conditions for guaranteeing the light-tail case in the solution scheme:

COROLLARY 1. *Suppose Assumptions 1 and 2 hold and (1) is consistent. An optimal solution for (1) must exist if $\lambda = 0$.*

Proof of Corollary 1. By Lemma 2, consistency of (1) implies $\sigma \geq \mu^2$. By Theorem 2 and Proposition 3, it suffices to consider the equivalent program $OPT(\mathcal{P}_2^+)$. Suppose $\lambda = 0$. If $\sigma = \mu^2$, then $\delta(\mu)$ is an optimal solution. If $\sigma > \mu^2$, then by Proposition 4, if there is no optimal solution, its optimal value must be $\nu(H(\mu) + \lambda(\sigma - \mu^2)) = \nu H(\mu)$, which is attained by $\delta(\mu)$ and leads to a contradiction (to both the hypotheses of no optimal solution and $\sigma > \mu^2$). □

We are now ready to show Theorem 4:

Proof of Theorem 4. Proof of 1. Assumption 1 implies Assumption 3. By Theorem 2 and Proposition 3, program (1) has the same optimal value as that of $OPT(\mathcal{P}_2^+)$. Similar to the proof of

Theorem 1, the result follows by noting that any $\mathbb{P} \in \mathcal{P}_2^+$ represented by (x_1, x_2, p_1, p_2) (with potentially identical x_i 's), admits one-to-one correspondence with a solution f in (1), via $f'_+(x+a) = \nu(p(x) - 1)$ in Theorem 2, giving

$$f'_+(x) = \begin{cases} -\nu & \text{for } a \leq x < x_1 + a \\ -\nu p_2 & \text{for } x_1 + a \leq x < x_2 + a \\ 0 & \text{for } x_2 + a \leq x \end{cases}$$

and hence

$$f(x) = \begin{cases} \eta - \nu(x - a) & \text{for } a \leq x \leq x_1 + a \\ \eta - \nu x_1 - \nu p_2(x - a - x_1) & \text{for } x_1 + a \leq x \leq x_2 + a \\ 0 & \text{for } x_2 + a \leq x \end{cases} \quad (10)$$

Proof of 2. The condition $\eta^2 < 2\beta\nu$ is equivalent to $\sigma > \mu^2$. The conclusion follows from Part 3 in Proposition 4.

Proof of 3. The first case is obtained by substituting $x_1^* \in \operatorname{argmax}_{x_1 \in [0, \mu]} W(x_1)$ and x_2^*, p_2^* from (9), in Part 1 in Proposition 4, into (10). The second case is obtained by substituting $(x_1^{(k)}, x_2^{(k)}, p_1^{(k)}, p_2^{(k)})$ in Part 2 in Proposition 4 into (10) and taking the limit. The last conclusion follows from Corollary 1. \square

6. Formulation and Procedure under Data-driven Environment

Sections 4 and 5 have discussed our worst-case approach in the abstract setting where the values of the needed parameters β, η, ν are completely known. In practice, these parameters are not directly specified. Instead, they are calibrated from data in the non-tail region. Suppose we obtain confidence intervals (CIs) for $P(X > a)$ and $f(a)$ and a lower confidence bound for $f'_-(a)$, jointly with confidence level $1 - \alpha$. Denote them as $[\underline{\beta}, \bar{\beta}]$, $[\underline{\eta}, \bar{\eta}]$ and $-\bar{\nu}$. Suppose $\underline{\beta}, \bar{\beta}, \underline{\eta}, \bar{\eta}, \bar{\nu} > 0$. We substitute

these estimates for the exact values of β , η and $-\nu$ in our worst-case bound for $E[h(X); X \geq a]$:

$$\begin{aligned}
& \max_f \quad \int_a^\infty h(x)f(x)dx \\
& \text{subject to} \quad \underline{\beta} \leq \int_a^\infty f(x)dx \leq \bar{\beta} \\
& \quad \underline{\eta} \leq f(a) = f(a+) \leq \bar{\eta} \\
& \quad f'_+(a) \geq -\bar{\nu} \\
& \quad f(x) \text{ convex for } x \geq a \\
& \quad f(x) \geq 0 \text{ for } x \geq a
\end{aligned} \tag{11}$$

It is immediate that the optimal value of (11) carries the following statistical guarantee:

PROPOSITION 5. *Suppose that $[\underline{\beta}, \bar{\beta}]$, $[\underline{\eta}, \bar{\eta}]$ and $-\bar{\nu}$ are the joint $(1 - \alpha)$ -level CIs for $P(X > a)$ and $f(a)$, and lower confidence bound for $f'_-(a)$. Then with probability $1 - \alpha$ (with respect to the data) optimization (11) gives an upper bound for $E[h(X); X \geq a]$ under the assumption that $f(x)$ is convex for $x \geq a$ and $f(a) = f(a+)$.*

Proof of Proposition 5. Let $f_{true}(x)$, $x \geq a$ be the ground-true density, and $Z_{true} = \int_a^\infty h(x)f_{true}(x)dx$. Let Z^* and \mathcal{F} be the optimal value and feasible region of (11). If $f_{true} \in \mathcal{F}$, then $Z^* \geq Z_{true}$. Hence $P_{\text{data}}(Z^* \geq Z_{true}) \geq P_{\text{data}}(f_{true} \in \mathcal{F}) = 1 - \alpha$, where P_{data} denotes the probability with respect to the data. \square

For h that has support spanning across both $X < a$ and $X \geq a$, one approach is to estimate $E[h(X); X < a]$ separately from the computation of the worst-case bound from (11). The former can be done typically by using the empirical mean as the non-tail region $X < a$ possesses more data to rely on. This segregated approach, however, only allows the conditions of valid probability density on the whole real line (e.g., $\int_{\mathbb{R}} f(x)dx = 1$) and the continuity at a to hold approximately but not exactly.

The following result presents the optimality structure for (11) in parallel to formulation (1).

THEOREM 5. *Suppose h is bounded. Consider optimization (11). If it is feasible, then either*

1. *An optimal solution f^* exists, where $f^* \in \mathcal{PL}_3^+[a, \infty)$.*

2. An optimal solution does not exist. There exists a sequence $\{f^{(k)} \in \mathcal{PL}_3^+[a, \infty) : k \geq 1\}$, each $f^{(k)}$ feasible for (1), such that $\int_a^\infty h(x)f^{(k)}(x)dx \rightarrow Z^*$ as $k \rightarrow \infty$, where Z^* is the optimal value of (11). Moreover, let $\{c_3^{(k)} + d_3^{(k)}x : x \in [y_2^{(k)}, y_3^{(k)}]\}$ be the last line segment of $f^{(k)}$. We have $y_3^{(k)} \nearrow \infty$ and $d_3^{(k)} \searrow 0$ as $k \rightarrow \infty$.

Proof of Theorem 5. See Appendix EC.3. □

Define

$$\underline{\mu} = \frac{\underline{\eta}}{\underline{\nu}}, \quad \bar{\mu} = \frac{\bar{\eta}}{\bar{\nu}}, \quad \underline{\sigma} = \frac{2\underline{\beta}}{\underline{\nu}}, \quad \bar{\sigma} = \frac{2\bar{\beta}}{\bar{\nu}} \quad (12)$$

where $\underline{\mu}, \bar{\mu}, \underline{\sigma}, \bar{\sigma} > 0$ since we have assumed $\underline{\beta}, \bar{\beta}, \underline{\eta}, \bar{\eta}, \underline{\nu} > 0$. Define

$$\mathcal{W}(x, \omega, \rho) = \bar{\nu} \left(\frac{\rho - \omega^2}{\rho - 2\omega x + x^2} H(x) + \frac{(\omega - x)^2}{\rho - 2\omega x + x^2} H\left(\frac{\rho - \omega x}{\omega - x}\right) \right)$$

with $\mathcal{W}(\omega, \omega, \rho) := \bar{\nu}(H(\omega) + \lambda(\rho - \omega^2))$, where H and λ are defined as in (3) and (7).

For convenience, we also denote

$$\mathcal{K}(x; x_1, \omega, \rho) = \begin{cases} \bar{\nu}\omega - \bar{\nu}(x - a) & \text{for } a \leq x \leq x_1 + a \\ \bar{\nu}\omega - \bar{\nu}x_1 - \bar{\nu} \frac{(\omega - x_1)^2}{\rho - 2\omega x_1 + x_1^2} (x - a - x_1) & \text{for } x_1 + a \leq x \leq \frac{\rho - \omega x_1}{\omega - x_1} + a \\ 0 & \text{for } \frac{\rho - \omega x_1}{\omega - x_1} + a \leq x \end{cases}$$

Our data-integrated optimization (11) possesses the following consistency property in parallel to the fixed-parameter case in Lemma 2:

LEMMA 3. Program (11) is consistent if and only if $\underline{\eta}^2 \leq 2\bar{\beta}\bar{\nu}$ or equivalently $\bar{\sigma} \geq \underline{\mu}^2$. When $\underline{\eta}^2 = 2\bar{\beta}\bar{\nu}$ or equivalently $\bar{\sigma} = \underline{\mu}^2$, (11) has only one feasible solution given by $f(x) = \underline{\eta} - \bar{\nu}(x - a)$ for $x \geq a$.

Proof of Lemma 3. The proof is similar to Lemma 2 and hence skipped. □

The following provides the solution scheme for our data-integrated optimization (11):

THEOREM 6. Under Assumption 1,

1. The conclusions of Theorem 5 hold with $\mathcal{PL}_3^+[a, \infty)$ replaced by $\mathcal{PL}_2^+[a, \infty)$.

2. Suppose $\underline{\eta}^2 < 2\bar{\beta}\bar{\nu}$ and Assumption 2 holds additionally. The optimal value of (11) is given by

$$\max \left\{ \max_{\rho \in [\underline{\sigma}\sqrt{\bar{\mu}^2}, \bar{\sigma}], x_1 \in [0, \bar{\mu}]} \mathcal{W}(x_1, \bar{\mu}, \rho), \max_{\omega \in [\underline{\mu}, \bar{\mu} \wedge \sqrt{\bar{\sigma}}], x_1 \in [0, \omega]} \mathcal{W}(x_1, \omega, \bar{\sigma}) \right\} \quad (13)$$

3. Suppose $\underline{\eta}^2 < 2\bar{\beta}\bar{\nu}$ and Assumption 2 holds additionally. Suppose $\max_{\rho \in [\underline{\sigma}\sqrt{\bar{\mu}^2}, \bar{\sigma}], x_1 \in [0, \bar{\mu}]} \mathcal{W}(x_1, \bar{\mu}, \rho) \geq \max_{\omega \in [\underline{\mu}, \bar{\mu} \wedge \sqrt{\bar{\sigma}}], x_1 \in [0, \omega]} \mathcal{W}(x_1, \omega, \bar{\sigma})$. If there exists $(\rho^*, x_1^*) \in \operatorname{argmax}_{\rho \in [\underline{\sigma}\sqrt{\bar{\mu}^2}, \bar{\sigma}], x_1 \in [0, \bar{\mu}]} \mathcal{W}(x_1, \bar{\mu}, \rho)$ such that $x_1^* \in [0, \bar{\mu})$, then an optimal solution to (11) is given by $f^*(x) = \mathcal{K}(x; x_1^*, \bar{\mu}, \rho^*)$. Otherwise, there exists a sequence of feasible solutions $f^{(k)}$ with $\int_a^\infty h(x)f^{(k)}(x)dx \rightarrow Z^*$, the optimal value of (11), such that $f^{(k)} \rightarrow f^*$ pointwise where

$$f^*(x) = \begin{cases} \bar{\eta} - \bar{\nu}(x - a) & \text{for } a \leq x \leq \bar{\mu} + a \\ 0 & \text{for } \bar{\mu} + a \leq x \end{cases}$$

which can occur only when $\lambda > 0$. On the other hand, suppose $\max_{\rho \in [\underline{\sigma}\sqrt{\bar{\mu}^2}, \bar{\sigma}], x_1 \in [0, \bar{\mu}]} \mathcal{W}(x_1, \bar{\mu}, \rho) < \max_{\omega \in [\underline{\mu}, \bar{\mu} \wedge \sqrt{\bar{\sigma}}], x_1 \in [0, \omega]} \mathcal{W}(x_1, \omega, \bar{\sigma})$. If there exists $(\omega^*, x_1^*) \in \operatorname{argmax}_{\omega \in [\underline{\mu}, \bar{\mu} \wedge \sqrt{\bar{\sigma}}], x_1 \in [0, \omega]} \mathcal{W}(x_1, \omega, \bar{\sigma})$ such that $x_1^* \in [0, \omega^*)$, then an optimal solution to (11) is given by $f^*(x) = \mathcal{K}(x; x_1^*, \omega^*, \bar{\sigma})$. Otherwise, there exists a sequence of feasible solutions $f^{(k)}$ with $\int_a^\infty h(x)f^{(k)}(x)dx \rightarrow Z^*$, such that $f^{(k)} \rightarrow f^*$ pointwise where

$$f^*(x) = \begin{cases} \bar{\nu}\omega^* - \bar{\nu}(x - a) & \text{for } a \leq x \leq \omega^* + a \\ 0 & \text{for } \omega^* + a \leq x \end{cases}$$

which again can occur only when $\lambda > 0$.

Proof of Theorem 6. Optimization (13) follows from a reduction of the inequality-based generalized moment problem converted from (11) into two subproblems. Appendix EC.3 provides the constituent propositions and further details. \square

Algorithm 2 presents our procedure for solving (11).

Algorithm 2: Procedure for Finding the Optimal Value of (11)

Inputs:

1. The function h that satisfies Assumptions 1 and 2.
2. The parameters $\underline{\beta}, \bar{\beta}, \underline{\eta}, \bar{\eta}, \bar{\nu} > 0$.

Procedure:

1. If $\underline{\eta}^2 > 2\bar{\beta}\bar{\nu}$, there is no feasible solution.
2. If $\underline{\eta}^2 = 2\bar{\beta}\bar{\nu}$, the optimal value is $\bar{\nu}H(\underline{\mu})$.
3. If $\underline{\eta}^2 < 2\bar{\beta}\bar{\nu}$, the optimal value is

$$\max \left\{ \max_{\rho \in [\underline{\sigma}\sqrt{\bar{\mu}^2}, \bar{\sigma}], x_1 \in [0, \bar{\mu}]} \mathcal{W}(x_1, \bar{\mu}, \rho), \max_{\omega \in [\underline{\mu}, \bar{\mu} \wedge \sqrt{\bar{\sigma}}], x_1 \in [0, \omega]} \mathcal{W}(x_1, \omega, \bar{\sigma}) \right\}$$

7. Numerical Examples

We present some numerical performance of our algorithm. We first consider several elementary examples, and then we will revisit the two examples in the Introduction. All computations are performed using the R software. For transparency and reproducibility purposes, the codes are made available on GitHub at the following url <https://github.com/cmottet/RobustTailNumerics>.

7.1. Elementary Examples

We consider three examples to demonstrate Algorithm 1.

Entropic Risk Measure: The entropic risk measure (e.g., Föllmer and Schied (2011)) captures the risk aversion of users through the exponential utility function. It is defined as

$$\rho(X) = \frac{1}{\theta} \log \left(E \left[e^{-\theta X} \right] \right) \quad (14)$$

where $\theta > 0$ is the parameter of risk aversion. In the case when the distribution of the random variable X is known only up to some point a , we can find the worst case value of the entropic risk measure subject to tail uncertainty by solving the optimization problem

$$\max_{P \in \mathcal{A}} \frac{1}{\theta} \log \left(E \left[e^{-\theta X} \right] \right) = \frac{1}{\theta} \log \left(E[e^{-\theta X}; X \leq a] + \max_{P \in \mathcal{A}} E[e^{-\theta X}; X > a] \right) \quad (15)$$

where \mathcal{A} denotes the set of convex tails that match the given non-tail region. Since the function $e^{-\theta X}$ satisfies Assumption 1, we can apply Algorithm 1 to the second term of the RHS of (15).

The thick line in Figure 9 represents the worst-case value of the entropic risk measure for different

values of the parameter θ in the case when X is known to have a standard exponential distribution $Exp(1)$ up to $a = -\log(0.7)$ (i.e. a is the 70-percentile and $\beta = \eta = \nu = 0.7$). For comparison, we also calculate and plot the entropic risk measure for several fitted probability distributions: $Exp(1)$, two-segment continuous piecewise linear tail denoted as 2-PLT (two such instances in Figure 9), and mixtures of 2-PLT and shifted Pareto. Clearly, the worst-case values bound those calculated from the candidate parametric models, with the gap diminishing as θ increases.

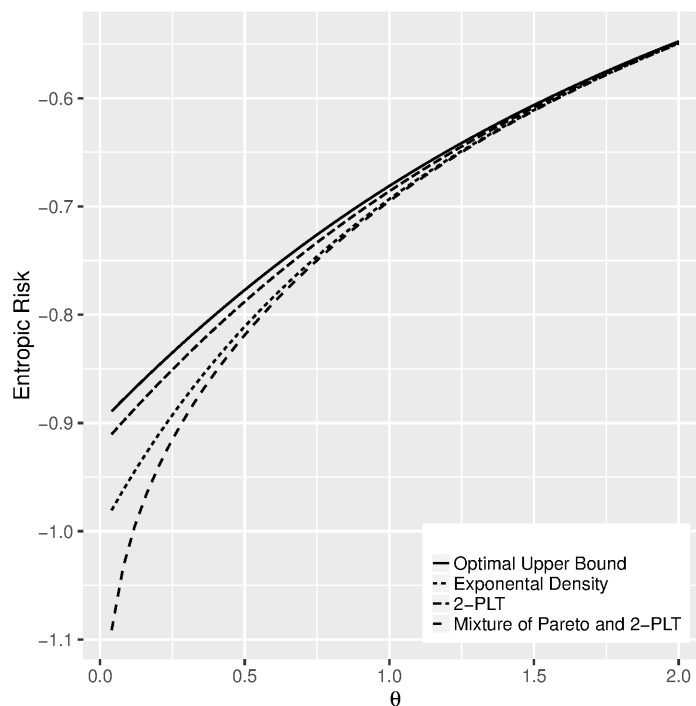


Figure 9 Optimal upper bound and comparison with parametric extrapolations for the entropic risk measure.

The Newsvendor Problem: The classical newsvendor problem maximizes the profit of selling a perishable product by fulfilling demand using a stock level decision, i.e.,

$$\max_q E[p \min(q, D)] - cq \quad (16)$$

where D is the demand random variable, p and c are the selling and purchase prices per product, and q is the stock quantity to be determined. We assume that $p > c$. The optimal solution to (16)

is given by Littlewood's rule $q^* = F^{-1}((p - c)/p)$, where F^{-1} is the quantile function of D (Talluri and Van Ryzin (2006)).

Suppose the distribution of D is only known to have the shape of a lognormal distribution with mean 50 and standard deviation 20 in the interval $[0, a]$, where a is the 70-percentile of the lognormal distribution. A robust optimization formulation for (16) is

$$\begin{aligned} & \max_q \min_{P \in \mathcal{A}} E[p \min(q, D)] - cq \\ & = \max_q \left\{ E[p \min(q, D); D \leq a] + \min_{P \in \mathcal{A}} E[p \min(q, D); D > a] - cq \right\} \end{aligned} \quad (17)$$

where \mathcal{A} denotes the set of convex tails that match the given non-tail region. The outer optimization in (17) is a concave program. We concentrate on the inner optimization. Since $p \min(q, D)$ is a non-decreasing function in D on $[0, \infty)$, its negation is non-increasing, and Assumption 1 holds (note that minimization here can be achieved by merely maximizing the negation). We can therefore apply Algorithm 1 (with $\beta = 0.7$, $\eta \approx 0.007$, and $\nu \approx 0.0003$). Figure 10 shows the optimal lower bound of the inner optimization when $p = 7$, $c = 1$ and q varies between 0 and 193.26 (which is the 95-percentile of the lognormal distribution). The curve peaks at $q = 55.7$, which is the solution to problem (17). As a comparison, we also show different candidate values of the expectation that are obtained by fitting the tails of lognormal, 2-PLT (two instances) and mixture of shifted Pareto and 2-PLT (see Figure 10).

Tail Interval Probability: Consider estimating probabilities of the type $P(c < X < d)$. We compare the bound provided by Algorithm 1 with the “truth” that X is a Pareto distribution with tail index 1, i.e. $P(X > x) = 1/x$ for all $x > 1$, for different values of the threshold a and interval (c, d) . On the heat map given in Figure 11, the color of each rectangle represents the ratio between the computed upper bound and the true probability for some threshold a and interval (c, d) . The x -value on the left hand side of each rectangle indicates the value of a (in percentile of the Pareto distribution) and the lower and upper y -values of the rectangle represent the interval (c, d) (also in percentile of the Pareto distribution). Each rectangle represents an area that has

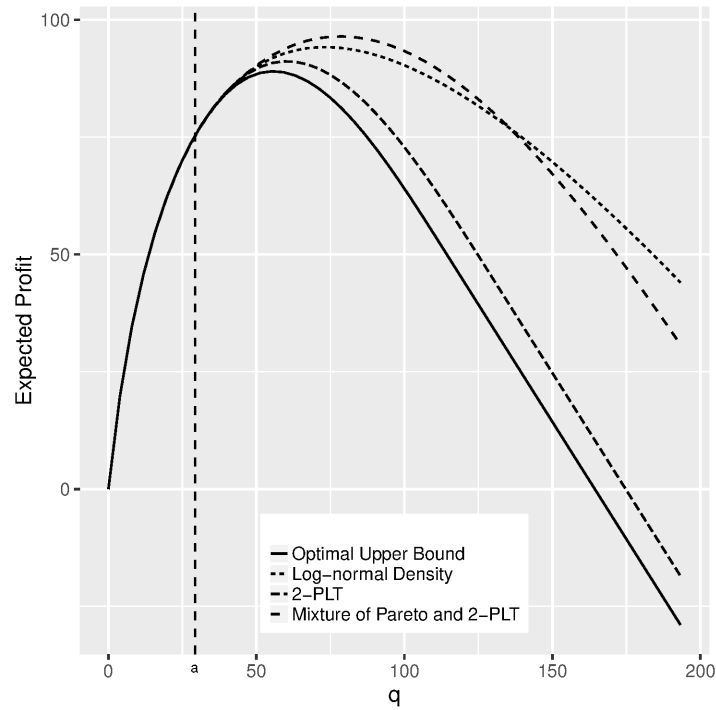


Figure 10 Optimal objective values of the inner optimization of the robust newsvendor problem.

Pareto probability mass exactly 0.01. We see that, for instance, when a is the 70th percentile and (c, d) is the (85th, 86th)-percentile, then our bound is roughly two times that of Pareto, whereas $(c, d) = (98^{\text{th}}, 99^{\text{th}})$ -percentile gives roughly eight times. Figure 11 confirms the intuition that the smaller the distance between a and c , the less conservative is the bound and hence the closer is to the “truth”.

7.2. Synthetic Data: Example 2 Revisited

Consider the synthetic data set of size 200 in Example 2. This data set is actually generated from a lognormal distribution with parameter $(\mu, \sigma) = (0, 0.5)$, but we assume that only the data are available to us. We are interested in the quantity $P(4 < X < 5)$, and for this we will solve program (11) to generate an upper bound that is valid with 95% confidence.

We compute the interval estimates for β , η and ν as follows. First, we obtain point estimates for these parameters through standard kernel density estimator (KDE) in the R statistical package. To obtain interval estimates, we run 1,000 bootstrap resamples and take the appropriate quantiles

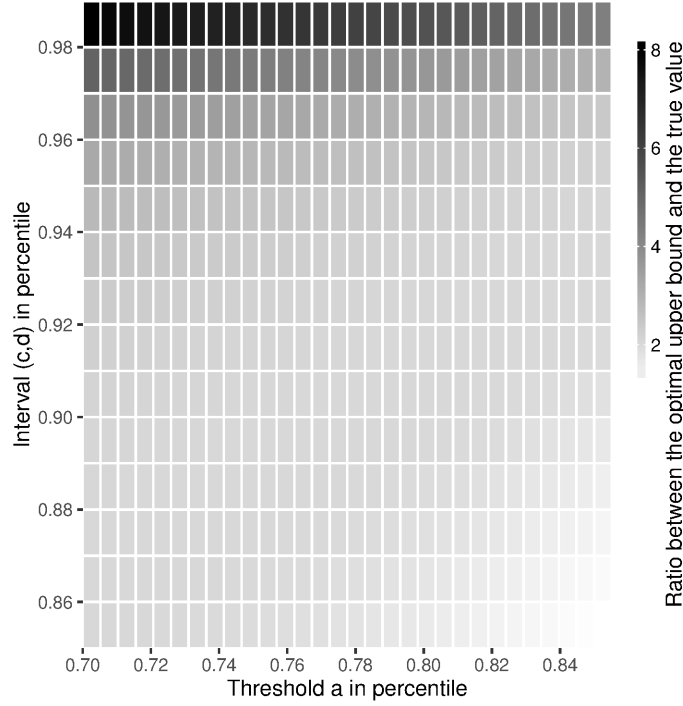


Figure 11 Ratio between the worst-case upper bound and the Pareto distribution for the quantity $P(c < X < d)$ for different thresholds a and intervals (c, d) .

of the 1,000 resampled point estimates. To account for the fact that three parameters are estimated simultaneously, we apply a Bonferroni correction, so that the confidence level used for each individual estimator is $1 - 0.05/3$.

For a sense of how to choose a , Figure 12 shows the density and density derivative estimates and compares them to those of the lognormal distribution. The KDE suggests that convexity holds starting from around $x = 1.5$ (the point where the density derivative estimate starts to turn from a decreasing to an increasing function). Thus, it is reasonable to confine the choice of a to be larger than 1.5. In fact, this number is quite close to the true inflexion point 1.15.

Since the data become progressively sparser as x grows larger, and the KDE is designed to utilize neighborhood data, the interval estimators for the necessary parameters β , η and ν become less reliable for larger choices of a . For instance, Figure 12 shows that the bootstrapped KDE CI of the density derivative covers the truth only up to $x = 3.1$. In general, a good choice of a should be located at a point where there are some data in the neighborhood of a , such that the interval

estimators for β , η and ν are reliable, but as large as possible, because choosing a small a can make the tail extrapolation bound more conservative.

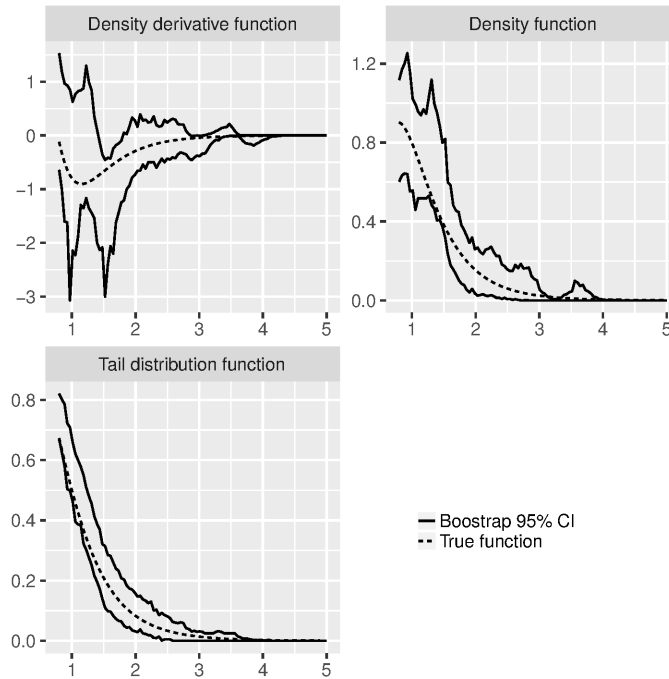


Figure 12 Bootstrapped kernel estimation of the distribution, density and density derivative for the synthetic data.

As a first attempt, we run Algorithm 2 using $a = 3.1$ to estimate an upper bound for the probability $P(4 < X < 5)$, which gives 8.8×10^{-3} while the truth is 2.1×10^{-3} . Thus, this estimated upper bound does cover the truth and also has the same order of magnitude. We perform the following two other procedures for comparison:

1. GPD approach: As discussed in Section 3.1, this is a common approach for tail modeling. Fit the data above a threshold u to the density function

$$(1 - \hat{F}(u))g_{\hat{\xi}, \hat{\beta}}(x - u)$$

where $\hat{F}(u)$ is the estimated ECDF at u , and $g_{\hat{\zeta},\hat{\beta}}(\cdot)$ is the GPD density, whose distribution function is defined as

$$G_{\zeta,\beta}(x) = \begin{cases} 1 - (1 + \zeta x/\beta)^{-1/\zeta} & \text{if } \zeta \neq 0 \\ 1 - \exp(-x/\beta) & \text{if } \zeta \geq 0 \end{cases}$$

for $x \geq 0$, and $\beta > 0$. Set the threshold u to be 1.8, the point at which a linear trend begins to be observed on the mean excess plot of the data, as recommended by McNeil (1997). Estimate $\hat{F}(u)$ by the sample mean of $I(X_i \leq u)$, where $I(\cdot)$ denotes the indicator function. Obtain the parameter estimates $\hat{\zeta}$ and $\hat{\beta}$ using the maximum likelihood estimator suggested by Smith (1987). Then use the delta method to obtain a 95% CI of the quantity $P(c < X < d)$.

2. Worst-case approach with known parameter values: Assume β , η and ν are known at $a = 3.1$. Then run Algorithm 1 to obtain the upper bound.

Table 1 shows the upper bounds obtained from the above approaches, and also shows the obvious fact that using ECDF alone for estimating $P(4 < X < 5)$ gives 0 since there are no data in the interval $[4, 5]$. The 95% CI output by GPD fit is $[-8.72 \times 10^{-4}, 1.10 \times 10^{-3}]$, which does not bound the truth (note that this is a two-sided interval, and the upper bound would be off even more if it had been one-sided). The worst-case approach with known parameters gives an upper bound of 3.16×10^{-3} , which is less conservative than the case when the parameters are estimated. The difference between these numbers can be interpreted as the price of estimation for β , η and ν . For this particular setup, the worst-case approach correctly covers the true value, whereas GPD fitting gives an invalid upper bound, thus showing that either the data size or the threshold level is insufficient to support a good fit of the GPD. This is an instance where the worst-case approach has outperformed GPD in terms of correctness.

Given that the worst-case approach with estimated parameters appears conceivably more conservative than with known parameters, we conduct a sensitivity study using only Algorithm 1. The first row in Table 2 shows the upper bound output by Algorithm 1 using the point estimates of the parameters β, η, ν . The other rows in Table 2 show the outputs of Algorithm 1 when some

Method	Estimated upper bound
Truth	2.14E-03
ECDF	0.00E+00
GPD	1.11E-03
Worst-case with known parameters	3.16E-03
Worst-case approach	8.80E-03

Table 1 Estimated upper bounds of the probability $P(4 < X < 5)$ for the synthetic data in Example 2.

values of the parameters are changed to the upper estimates of the 95% CIs. Some scenarios are omitted in the table because they lead to infeasibility. We see that among all these scenarios, the most conservative upper bound occurs when β, η, ν are all set to be the upper estimates, giving to 8.67×10^{-3} which is very close to using Algorithm 2. Note that some of these bounds do not cover the truth, which necessitates the use of the interval approach and Algorithm 2.

The above discussion focuses only on the realization of one data set, which raises the question of whether it holds more generally. Therefore, we obtain an empirical probability of coverage by repeating the following procedure 100 times:

1. Generate a lognormal sample of size 200 with parameters $(\mu, \sigma) = (0, 0.5)$;
2. Estimate $\bar{\eta}$, $\underline{\eta}$, $\bar{\beta}$, $\underline{\beta}$ and $\bar{\nu}$ at a chosen point a (see below);
3. Use Algorithm 2 to compute the worst-case upper bound of $P(c < X < d)$.

We then estimate the coverage probability of our worst-case upper bound as the proportion of times that Algorithm 2 yields a bound that dominates the true probability $P(c < X < d)$. We repeat this procedure for different $[c, d]$ varying from $[4, 5]$ to $[9, 10]$, and for two different values of a given by 3.1 and 2.8. Tables 3 and 4 show the true probabilities, the mean upper bounds from the 100 experiments, and the empirical coverage probabilities.

The coverage probabilities in Tables 3 and 4 are mostly 1, which suggests that our procedure is conservative. For $a = 3.1$ and intervals that are close to a , i.e. $[c, d] = [4, 5]$ and $[5, 6]$, the coverage probability is not 1 but rather is close to the prescribed confidence level of 95%. Further investigation reveals that our procedure fails to cover the truth only in the case when the joint CI of the

β	η	ν	Worst-case upper bound
Estimated value	Estimated value	Estimated value	2.04E-03
Estimated value	Lower estimate	Estimated value	5.76E-06
Estimated value	Lower estimate	Upper estimate	5.76E-06
Upper estimate	Lower estimate	Estimated value	5.76E-06
Upper estimate	Lower estimate	Upper estimate	5.76E-06
Estimated value	Upper estimate	Estimated value	3.61E-04
Estimated value	Upper estimate	Upper estimate	1.62E-03
Estimated value	Estimated value	Upper estimate	2.05E-03
Upper estimate	Estimated value	Upper estimate	5.53E-03
Upper estimate	Estimated value	Estimated value	5.53E-03
Upper estimate	Upper estimate	Estimated value	8.30E-03
Upper estimate	Upper estimate	Upper estimate	8.67E-03

Table 2 Sensitivity analysis of the worst-case upper bound of $P(4 < X < 5)$ for the synthetic data in Example 2 generated by Algorithm 1, when β, η, ν are changed from the point estimates to the upper estimates of the 95% CIs.

c	d	Truth	Mean upper bound	Coverage probability
4	5	2.14E-03	1.03E-02	0.94
5	6	4.74E-04	6.12E-03	0.99
6	7	1.20E-04	4.33E-03	1.00
7	8	3.38E-05	3.35E-03	1.00
8	9	1.04E-05	2.74E-03	1.00
9	10	3.49E-06	2.31E-03	1.00

Table 3 Mean upper bounds and empirical coverage probabilities using worst-case approach with threshold

$$a = 3.1.$$

c	d	Truth	Mean upper bound	Coverage probability
4	5	2.14E-03	1.31E-02	1.00
5	6	4.74E-04	8.26E-03	1.00
6	7	1.20E-04	6.04E-03	1.00
7	8	3.38E-05	4.76E-03	1.00
8	9	1.04E-05	3.92E-03	1.00
9	10	3.49E-06	3.34E-03	1.00

Table 4 Mean upper bounds and empirical coverage probabilities using worst-case approach with threshold $a = 2.8$.

parameters η , β and ν does not contain the true values, which is consistent with the rationale of our method. Although we have not tried lower values of a , it is very likely that in those settings the coverage probabilities will stay mostly 1, and the mean upper bounds will increase since the level of conservativeness increases.

As a comparison, Table 5 shows the results of GPD fit using the threshold $u = 1.8$. Here, all of the coverage probabilities are far from the prescribed level of 95%, which suggests that either GPD is the wrong parametric choice to use since the threshold is not high enough, or that the estimation error of its parameters is too large due to the lack of data. (Again, we have used a two-sided 95% CI for the GPD approach here; if we had used a one-sided upper confidence bound, then the upper bounding value would be even lower and the coverage probability would drop further). However, the mean upper bounds using GPD fit do cover the truth in all cases. Since the coverage probability is well below 95%, this suggests that the estimation of GPD parameters is highly sensitive to the realization of data.

In summary, Tables 3, 4 and 5 show the pros and cons of our worst-case approach and GPD fitting. GPD is on average closer to the true target quantity, but its confidence upper bound can fall short of the prescribed coverage probability (in fact, only between 37 to 62% of the time it covers the truth in Table 5). On the other hand, our approach gives a reasonably tight upper bound when

c	d	Truth	Mean upper bound	Coverage probability
4	5	2.14E-03	3.87E-03	0.62
5	6	4.74E-04	1.27E-03	0.53
6	7	1.20E-04	5.48E-04	0.51
7	8	3.38E-05	2.79E-04	0.43
8	9	1.04E-05	1.62E-04	0.40
9	10	3.49E-06	1.03E-04	0.37

Table 5 Mean upper bounds and empirical coverage probabilities using GPD.

the interval in consideration (i.e. $[c, d]$) is close to the threshold a , and tends to be more conservative far out. This is a drawback, but sensibly so, given that the uncertainty of extrapolation increases as it gets farther away from what is known.

Both our worst-case approach and GPD fitting require choosing a threshold parameter. In GPD fitting, it is important to choose a threshold parameter high enough so that the GPD becomes a valid model. GPD fitting, however, is difficult for a small data set when the lack of data prohibits choosing a high threshold. On the other hand, the threshold in our worst-case approach can be chosen much higher, because our method relies on the data below the threshold, not above it.

7.3. Fire Insurance Data: Example 1 Revisited

Consider the fire insurance data in Example 1. The quantity of interest is the expected payoff of a high-excess policy with reinsurance, given by $h(x) = (x - 50)I(50 \leq x < 200) + 150I(x \geq 200)$. The data set has only seven observations above 50.

We apply our worst-case approach to estimate an upper bound for the expected payoff by using $a = 29.03$, the cutoff above which 15 observations are available. Similar to Section 7.2, we use the bootstrapped KDE to obtain CIs for β , η and ν . The estimates in Figure 13 appear to be very stable for this example, thanks to the relatively large data size.

We run Algorithm 2 to obtain a 95% confidence upper bound of 1.99. For comparison, we fit a GPD using threshold $u = 10$, which follows McNeil (1997) as the choice that roughly balances the

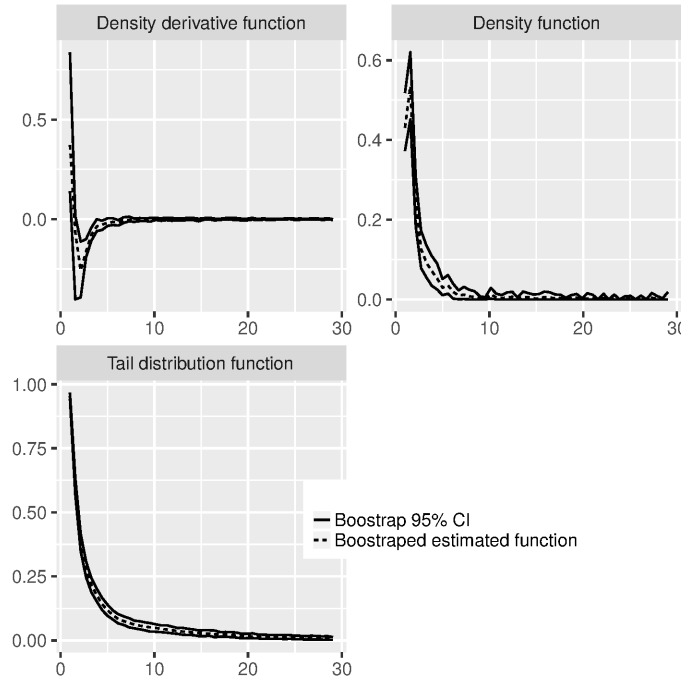


Figure 13 Bootstrapped kernel estimation of the distribution, density and density derivative for the the Danish fire losses data in Example 1.

bias-variance tradeoff. The 95% CI from GPD fit is $[-0.03, 0.23]$. Thus, the worst-case approach gives an upper bound that is one order of magnitude higher, a finding that resonates with that in Section 7.2. Our recommendation is that a modeler who cares only about the order of magnitude would be better off choosing GPD, whereas a more risk-averse modeler who wants a bound on the risk quantity with high probability guarantee would be better off choosing the worst-case approach.

8. Conclusion

This paper proposed a worst-case, nonparametric approach to bound tail quantities based on the tail convexity assumption. The approach relied on an optimization formulated over all possible tail densities. We characterized the optimality structure of this infinite-dimensional optimization problem by developing an equivalence to a moment-constrained problem. Under additional quasi-concavity condition on the objective function, we constructed the numerical solution scheme by converting it into low-dimensional nonlinear programs. With the presence of data, this approach tractably generated statistically valid bounds via suitable relaxations of the optimization that took

into account the estimation errors of the required parameters. We compared our proposed approach to existing tail-fitting techniques, and demonstrated its relative strength of outputting correct tail estimates under data-deficient environments. We also examined the level of conservativeness of our bounds, which was viewed as a limitation of the proposed approach.

We suggest two extensions of our research. First is to generalize the proposed method to multivariate distributions, perhaps through separate modeling on the marginal distributions and the dependency structure. Second is to study means to reduce the level of conservativeness. This can involve mathematical transformations of the variable and the addition of extra information (e.g., other constraints).

References

- Balkema, August A, Laurens De Haan. 1974. Residual life time at great age. *The Annals of Probability* 792–804.
- Beirlant, Jan, Jozef L Teugels. 1992. Modeling large claims in non-life insurance. *Insurance: Mathematics and Economics* **11**(1) 17–29.
- Ben-Tal, Aharon, Dick Den Hertog, Anja De Waegenaere, Bertrand Melenberg, Gijs Rennen. 2013. Robust solutions of optimization problems affected by uncertain probabilities. *Management Science* **59**(2) 341–357.
- Bertsimas, Dimitris, Ioana Popescu. 2005. Optimal inequalities in probability theory: A convex optimization approach. *SIAM Journal on Optimization* **15**(3) 780–804.
- Birge, John R, José H Dulá. 1991. Bounding separable recourse functions with limited distribution information. *Annals of Operations Research* **30**(1) 277–298.
- Birge, John R, Roger J-B Wets. 1987. Computing bounds for stochastic programming problems by means of a generalized moment problem. *Mathematics of Operations Research* **12**(1) 149–162.
- Cule, Madeleine, Richard Samworth, Michael Stewart. 2010. Maximum likelihood estimation of a multi-dimensional log-concave density. *Journal of the Royal Statistical Society: Series B (Statistical Methodology)* **72**(5) 545–607.

- Davis, Richard, Sidney Resnick. 1984. Tail estimates motivated by extreme value theory. *The Annals of Statistics* 1467–1487.
- Davison, Anthony C, Richard L Smith. 1990. Models for exceedances over high thresholds. *Journal of the Royal Statistical Society. Series B (Methodological)* 393–442.
- Delage, Erick, Yinyu Ye. 2010. Distributionally robust optimization under moment uncertainty with application to data-driven problems. *Operations Research* **58**(3) 595–612.
- El Ghaoui, Laurent, A Nilim. 2005. Robust solutions to Markov decision problems with uncertain transition matrices. *Operations Research* **53**(5).
- Embrechts, Paul, Rdiger Frey, Alexander McNeil. 2005. Quantitative risk management. *Princeton Series in Finance, Princeton* **10**.
- Embrechts, Paul, Claudia Klüppelberg, Thomas Mikosch. 1997. *Modelling extremal events*, vol. 33. Springer Science & Business Media.
- Fisher, Ronald Aylmer, Leonard Henry Caleb Tippett. 1928. Limiting forms of the frequency distribution of the largest or smallest member of a sample. *Mathematical Proceedings of the Cambridge Philosophical Society*, vol. 24. Cambridge University Press, 180–190.
- Föllmer, Hans, Alexander Schied. 2011. *Stochastic finance: an introduction in discrete time*. Walter de Gruyter.
- Glasserman, Paul, Wanmo Kang, Perwez Shahabuddin. 2007. Large deviations in multifactor portfolio credit risk. *Mathematical Finance* **17**(3) 345–379.
- Glasserman, Paul, Wanmo Kang, Perwez Shahabuddin. 2008. Fast simulation of multifactor portfolio credit risk. *Operations Research* **56**(5) 1200–1217.
- Glasserman, Paul, Jingyi Li. 2005. Importance sampling for portfolio credit risk. *Management Science* **51**(11) 1643–1656.
- Gnedenko, Boris. 1943. Sur la distribution limite du terme maximum d’une serie aleatoire. *Annals of Mathematics* 423–453.
- Goh, Joel, Melvyn Sim. 2010. Distributionally robust optimization and its tractable approximations. *Operations Research* **58**(4-part-1) 902–917.

- Gumbel, Emil Julius. 2012. *Statistics of Extremes*. Courier Corporation.
- Hannah, Lauren A, David B Dunson. 2013. Multivariate convex regression with adaptive partitioning. *The Journal of Machine Learning Research* **14**(1) 3261–3294.
- Hansen, Lars Peter, Thomas J Sargent. 2008. *Robustness*. Princeton University Press.
- Heidelberger, Philip. 1995. Fast simulation of rare events in queueing and reliability models. *ACM Transactions on Modeling and Computer Simulation (TOMACS)* **5**(1) 43–85.
- Hill, Bruce M, et al. 1975. A simple general approach to inference about the tail of a distribution. *The Annals of Statistics* **3**(5) 1163–1174.
- Hogg, Robert V, Stuart A Klugman. 2009. *Loss Distributions*, vol. 249. John Wiley & Sons.
- Hosking, Jonathan RM, James R Wallis. 1987. Parameter and quantile estimation for the generalized pareto distribution. *Technometrics* **29**(3) 339–349.
- Iyengar, Garud N. 2005. Robust dynamic programming. *Mathematics of Operations Research* **30**(2) 257–280.
- Karr, Alan F. 1983. Extreme points of certain sets of probability measures, with applications. *Mathematics of Operations Research* **8**(1) 74–85.
- Koenker, Roger, Ivan Mizera. 2010. Quasi-concave density estimation. *The Annals of Statistics* 2998–3027.
- Lim, Eunji, Peter W Glynn. 2012. Consistency of multidimensional convex regression. *Operations Research* **60**(1) 196–208.
- McNeil, A. J. 1997. Estimating the tails of loss severity distributions using extreme value theory. *The Journal of the International Actuarial Association* **27** 117–137.
- Nicola, Victor F, Marvin K Nakayama, Philip Heidelberger, Ambuj Goyal. 1993. Fast simulation of highly dependable systems with general failure and repair processes. *IEEE Transactions on Computers* **42**(12) 1440–1452.
- Petersen, Ian R, Matthew R James, Paul Dupuis. 2000. Minimax optimal control of stochastic uncertain systems with relative entropy constraints. *IEEE Transactions on Automatic Control* **45**(3) 398–412.
- Pickands III, James. 1975. Statistical inference using extreme order statistics. *The Annals of Statistics* 119–131.

- Popescu, Ioana. 2005. A semidefinite programming approach to optimal-moment bounds for convex classes of distributions. *Mathematics of Operations Research* **30**(3) 632–657.
- Rockafellar, Ralph Tyrell. 2015. *Convex analysis*. Princeton university press.
- Seijo, Emilio, Bodhisattva Sen, et al. 2011. Nonparametric least squares estimation of a multivariate convex regression function. *The Annals of Statistics* **39**(3) 1633–1657.
- Seregin, Arseni, Jon A Wellner. 2010. Nonparametric estimation of multivariate convex-transformed densities. *Annals of statistics* **38**(6) 3751.
- Smith, James E. 1995. Generalized chebychev inequalities: theory and applications in decision analysis. *Operations Research* **43**(5) 807–825.
- Smith, Richard L. 1985. Maximum likelihood estimation in a class of nonregular cases. *Biometrika* **72**(1) 67–90.
- Smith, Richard L. 1987. Estimating tails of probability distributions. *The annals of Statistics* 1174–1207.
- Talluri, Kalyan T, Garrett J Van Ryzin. 2006. *The theory and practice of revenue management*, vol. 68. Springer Science & Business Media.
- Winkler, Gerhard. 1988. Extreme points of moment sets. *Mathematics of Operations Research* **13**(4) 581–587.

Appendix

EC.1. Proofs for Section 4

We need several results from convex analysis to prove Lemma 1. For any convex function g on \mathbb{R} , let $\text{dom } g = \{x \in \mathbb{R} : g(x) < \infty\}$ be its effective domain. The following theorems are from Rockafellar (2015), specialized to convex functions g with $\text{dom } g = \mathbb{R}$.

THEOREM EC.1 (a.k.a. Rockafellar (2015), Corollary 10.1.1). *A convex function finite on \mathbb{R} is necessarily continuous.*

THEOREM EC.2 (a.k.a. Rockafellar (2015), Theorem 24.1). *Let g be a closed proper convex function on \mathbb{R} , such that $\text{dom } g = \mathbb{R}$. Then g'_+ exists and is a finite non-decreasing function on \mathbb{R} . Moreover, g'_+ is right-continuous, i.e., $\lim_{z \searrow x} g'_+(z) = g'_+(x)$ for any $x \in \mathbb{R}$.*

THEOREM EC.3 (a.k.a. Rockafellar (2015), Corollary 24.2.1). *Let g be a finite convex function on a non-empty open real interval I . Then*

$$g(y) - g(x) = \int_x^y g'_+(t) dt$$

for any x and y in I .

THEOREM EC.4 (a.k.a. Rockafellar (2015), Theorem 24.2). *Let φ be a non-decreasing function from \mathbb{R} to $[-\infty, \infty]$ such that $\varphi(b)$ is finite for some $b \in \mathbb{R}$. Then the function given by*

$$g(x) = \int_b^x \varphi(t) dt$$

is a well-defined closed proper convex function on \mathbb{R} .

Proof of Lemma 1. Throughout this proof, without loss of generality let $a = 0$ (by replacing $f(x)$ with $f(x + a)$, and $h(x)$ with $h(x + a)$ respectively). Note that optimizations (1) and (2) do not depend on $f(x)$ for $x < 0$. For the purpose of applying Theorems EC.1–EC.4 more directly, let us extend f to \mathbb{R}^- , by defining $f(x) = \eta - \nu x$ for $x < 0$ (this extension of f is a mathematical artifact and does not necessarily match the given true density).

Let \mathcal{F}_1 be the feasible region in (1), and \mathcal{F}_2 be the feasible region in (2). We show that $\mathcal{F}_1 = \mathcal{F}_2$.

Proof of $\mathcal{F}_1 \subset \mathcal{F}_2$: Since $f(x) < \infty$ for at least one $x \in \mathbb{R}$ (e.g., take $x = 0$) and $f(x) \geq 0 > -\infty$ for all $x \in \mathbb{R}$, we get that f is proper (Rockafellar (2015), p.24).

Next, we argue that $f(x) < \infty$ for all $x \in \mathbb{R}$. Suppose on the contrary that $f(x_0) = \infty$ for some $x > 0$. If $f(y) < \infty$ for some $y > x_0$, then $((y - x_0)/y)f(0) + (x_0/y)f(y) = ((y - x_0)/y)\eta + (x_0/y)f(y) < \infty = f(x_0)$, contradicting (1d). But if $f(y) = \infty$ for all $y > x_0$, then $\int_0^\infty f(t)dt = \infty$, contradicting (1a). Therefore, $f(x) < \infty$ for all $x \in \mathbb{R}$, and with (1e), we conclude that f is finite.

Since f is proper, closedness is the same as lower semi-continuity (Rockafellar (2015), p.52). Since f is finite on \mathbb{R} , Theorem EC.1 implies that f is continuous. Hence f is closed.

Therefore, together with the convexity condition in (1d), Theorem EC.2 implies the existence of f'_+ that satisfies (2c). Moreover, Theorem EC.3 implies (2f).

Next, with the monotonicity of f'_+ by (2c), we have $f'_+(x) \geq f'_+(0) = -\nu$ for all $x \geq 0$, thus implying the first inequality of (2d). To prove the second inequality in (2d), suppose in the contrary that $f'_+(x_0) > 0$ for some $x_0 > 0$. Since $f'_+(x) \geq f'_+(x_0) > 0$ for all $x > x_0$ by (2c), we have, from (2f), $f(x) = \int_0^x f'_+(t)dt + \eta \rightarrow \infty$, implying that $\int_0^\infty f(x)dx = \infty$ and contradicting (1a). Hence the second inequality in (2d) holds. We have therefore shown (2d).

Lastly, suppose that $f'_+(x) \not\rightarrow 0$. Then, since (2d) holds, there exists a sequence $x_k \rightarrow \infty$ such that $f'_+(x_k) \rightarrow c < 0$. But since f'_+ is monotone by (2c), $\lim_{x \rightarrow \infty} f'_+(x)$ exists and must equal c . But then, from (2f), $f(x) = \int_0^x f'_+(t)dt + \eta \rightarrow -\infty$, violating (1e). Thus (2e) holds.

The constraints (2a) and (2b) follow immediately from (1a) and (1b). We therefore conclude that $\mathcal{F}_1 \subset \mathcal{F}_2$.

Proof of $\mathcal{F}_2 \subset \mathcal{F}_1$: Since f'_+ is bounded on \mathbb{R} by (2d), Theorem EC.4 and (2c) (with $f'_+(x)$ defined as $-\nu$ for $x < 0$) implies that the f defined by (2f) is convex on \mathbb{R} , giving (1d).

Suppose $f(x_0) < 0$ for some $x_0 > 0$. Then, since $f'_+ \leq 0$ by (2d), (2f) implies $f(x) < 0$ for all $x \geq x_0$. Thus $\int_0^\infty f(x)dx = -\infty$, contradicting (2a). Therefore, (1e) holds.

The constraint (2d) implies (1c) immediately. The condition (2f) implies $f(0) = f(0+)$. Thus, combining with (2b), we get that (1b) holds. Finally, note that (2a) is the same as (1a). We conclude that $\mathcal{F}_2 \subset \mathcal{F}_1$. □

Proof of Theorem 2. Throughout this proof, without loss of generality let $a = 0$. For convenience, we let $\tilde{H}(x) = \int_0^x h(u)du$ and $H(x) = \int_0^x \tilde{H}(u)du$. Consider the objective function of (2). Since \tilde{H} is continuous and f is absolutely continuous with $f(x) = \int_0^x f'_+(t)dt + \eta$ by (2f), we have, using integration by parts,

$$\int_0^\infty h(x)f(x)dx = \tilde{H}(x)f(x)\Big|_0^\infty - \int_0^\infty \tilde{H}(x)f'_+(x)dx = - \int_0^\infty \tilde{H}(x)f'_+(x)dx \quad (\text{EC.1})$$

where the second equality follows from Lemma EC.1 (presented next) and that $\tilde{H}(x) = O(x)$ as $x \rightarrow \infty$ since h is bounded. As H is continuous and f'_+ has bounded variation by (2d) and (2c), we have, using integration by parts again, that (EC.1) is equal to

$$-H(x)f'_+(x)\Big|_0^\infty + \int_0^\infty H(x)df'_+(x) = \int_0^\infty H(x)df'_+(x) \quad (\text{EC.2})$$

where the equality follows from Lemma EC.1 and that $H(x) = O(x^2)$ as $x \rightarrow \infty$ since h is bounded.

For (2a), we can write

$$\int_0^\infty f(x)dx = \int_0^\infty \frac{x^2}{2}df'_+(x) \quad (\text{EC.3})$$

by merely viewing $h \equiv 1$ in (EC.1) and (EC.2). Also, since $f(x) \rightarrow 0$ as $x \rightarrow \infty$ by Lemma EC.1, we can use integration by parts again to write

$$f(0) = - \int_0^\infty f'_+(x)dx = -xf'_+(x)\Big|_0^\infty + \int_0^\infty xdf'_+(x) = \int_0^\infty xdf'_+(x) \quad (\text{EC.4})$$

where the third equality follows from Lemma EC.1 again. Therefore, (2) can be written as (letting $a = 0$)

$$\begin{aligned} & \max_f && \int_0^\infty H(x)df'_+(x) \\ & \text{subject to} && \int_0^\infty \frac{x^2}{2}df'_+(x) = \beta \\ & && \int_0^\infty xdf'_+(x) = \eta \\ & && f'_+(x) \text{ exists and is non-decreasing and right-continuous for } x \geq 0 \\ & && -\nu \leq f'_+(x) \leq 0 \text{ for } x \geq 0 \\ & && f'_+(x) \rightarrow 0 \text{ as } x \rightarrow \infty \end{aligned} \quad (\text{EC.5})$$

and the last constraint (2f) in (2) states that f can be recovered from $f(x) = \int_0^x f'_+(t)dt + \eta$. Note that this definition of f must necessarily have a right derivative coinciding with the obtained $f'_+(x)$.

Finally, let $p(x) = f'_+(x)/\nu + 1$. Then (EC.5) can be rewritten as

$$\begin{aligned}
& \max_p && \nu \int_0^\infty H(x) dp(x) \\
& \text{subject to} && \int_0^\infty x^2 dp(x) = \frac{2\beta}{\nu} \\
& && \int_0^\infty x dp(x) = \frac{\eta}{\nu} \\
& && p(x) \text{ non-decreasing and right-continuous for } x \geq 0 \\
& && 0 \leq p(x) \leq 1 \text{ for } x \geq 0 \\
& && p(x) \rightarrow 1 \text{ as } x \rightarrow \infty
\end{aligned} \tag{EC.6}$$

or equivalently

$$\begin{aligned}
& \max_p && \nu \int_{-\infty}^\infty H(x) dp(x) \\
& \text{subject to} && \int_{-\infty}^\infty x^2 dp(x) = \frac{2\beta}{\nu} \\
& && \int_{-\infty}^\infty x dp(x) = \frac{\eta}{\nu} \\
& && p(x) \text{ non-decreasing and right-continuous for } x \in \mathbb{R} \\
& && 0 \leq p(x) \leq 1 \text{ for } x \in \mathbb{R} \\
& && p(x) \rightarrow 1 \text{ as } x \rightarrow \infty \\
& && p(x) = 0 \text{ for } x < 0
\end{aligned} \tag{EC.7}$$

since $H(x) = x = x^2 = 0$ at $x = 0$. One can uniquely identify, up to measure zero, a non-decreasing, right-continuous p such that $\lim_{x \rightarrow \infty} p(x) = 1$ and $p(x) = 0$ for $x < 0$ with a probability measure supported on \mathbb{R}^+ . Hence (EC.7) is equivalent to (5). This concludes the result. \square

LEMMA EC.1. *If f is a feasible solution of (1), equivalently (2), then $xf(x) \rightarrow 0$ and $x^2 f'_+(x) \rightarrow 0$ as $x \rightarrow \infty$.*

Proof of Lemma EC.1. We need the observations that $f(x)$ is non-increasing by (2d), $f(x) \geq 0$ for all $x \geq a$ by (1e), and that f is integrable on $[a, \infty)$ with $\int_a^\infty f(x)dx = \beta$ by (1a). Denote $F(x) = \int_a^x f(t)dt$ and $g(x) = xf(x) - F(x)$. Consider, for $a \vee 0 \leq x_1 \leq x_2$,

$$g(x_2) - g(x_1) = x_2 f(x_2) - x_1 f(x_1) - (F(x_2) - F(x_1))$$

$$\begin{aligned}
&\leq x_2 f(x_2) - x_1 f(x_1) - f(x_2)(x_2 - x_1) \quad \text{since } f(x) \text{ is non-increasing} \\
&= x_1 [f(x_2) - f(x_1)] \\
&\leq 0 \quad \text{again since } f \text{ is non-increasing}
\end{aligned}$$

Therefore g is non-increasing for $x \geq a \vee 0$, and since $xf(x) \geq 0$ and $0 \leq F(x) \leq \beta$ for $x \geq a \vee 0$, we have g bounded from below on the same range. This implies that g must converge to a limit, say c , as $x \rightarrow \infty$. In other words, $xf(x) - F(x) \rightarrow c$, and since $F(x) \rightarrow \beta$, we have $xf(x) \rightarrow c + \beta$. Since $xf(x) \geq 0$ for $x \geq a \vee 0$, there are two cases: $c + \beta > 0$ or $c + \beta = 0$. The first case implies that $xf(x) \geq \epsilon > 0$ for some ϵ for all large enough x . This means $f(x) \geq \epsilon/x$ for all large enough x , and hence $\int_a^\infty f(x)dx = \infty$, which contradicts (1a). Therefore $xf(x)$ must converge to 0. This proves the first part of the lemma.

To prove the second part, we need the observation that $f'_+(x)$ is non-decreasing for $x \geq a$ by (2c), and is non-positive for $x \geq a$ by (2d). Also, by (2f) we have $f(x) = \int_a^x f'_+(t)dt + \eta$ for $x \geq a$. Let $\bar{F}(x) = \int_x^\infty f(t)dt$ for $x \geq a$, which is finite and converges to 0 by (1a). We now define $\tilde{g}(x) = -x^2 f'_+(x) + 2\tilde{F}(x)$, where $\tilde{F}(x) = -\int_x^\infty t f'_+(t)dt$, for $x \geq a$. Note that $xf'_+(x)$ is integrable on $[a, \infty)$ because the absolute continuity of f , and $\lim_{x \rightarrow \infty} xf(x) \rightarrow 0$ as we have just proved, which allows integration by parts yielding

$$\tilde{F}(x) = -\int_x^\infty t f'_+(t)dt = -tf(t)|_x^\infty + \int_x^\infty f(t)dt = xf(x) + \bar{F}(x) < \infty \quad (\text{EC.8})$$

For any $(a \vee 0) \leq x_1 \leq x_2$,

$$\begin{aligned}
\tilde{g}(x_2) - \tilde{g}(x_1) &= x_1^2 f'_+(x_1) - x_2^2 f'_+(x_2) - 2\tilde{F}(x_1) + 2\tilde{F}(x_2) \\
&\leq x_1^2 f'_+(x_1) - x_2^2 f'_+(x_2) + f'_+(x_2)(x_2^2 - x_1^2) \quad \text{since } f'_+(x) \text{ is non-decreasing} \\
&= x_1^2 (f'_+(x_1) - f'_+(x_2)) \\
&\leq 0 \quad \text{again since } f'_+(x) \text{ is non-decreasing}
\end{aligned}$$

Therefore, $\tilde{g}(x)$ is non-increasing for $x \geq a$. Note that $-x^2 f'_+(x) \geq 0$ for $x \geq a$. Also, from (EC.8), since $\lim_{x \rightarrow \infty} xf(x) \rightarrow 0$ and $\bar{F}(x) \rightarrow 0$, we have $\tilde{F}(x) \rightarrow 0$ as $x \rightarrow \infty$ and hence also bounded for

large enough x . Therefore \tilde{g} is bounded from below. This implies that \tilde{g} must converge to a limit, say \tilde{c} , as $x \rightarrow \infty$. Since $\tilde{F}(x) \rightarrow 0$, we have $-x^2 f'_+(x) \rightarrow \tilde{c}$. Since $-x^2 f'_+(x) \geq 0$ for $x \geq a$, there are two cases: either $\tilde{c} > 0$ or $\tilde{c} = 0$. The former case implies that $-x f'_+(x) \geq \epsilon/x$ for some $\epsilon > 0$ and large enough x , and so $\tilde{F}(x) = -\int_x^\infty x f'_+(x) dx = \infty$ for $x \geq a$, which contradicts (EC.8). Therefore $-x^2 f'_+(x) \rightarrow 0$. This proves the second part of the lemma. \square

To prove Theorem 3, we need several results from Winkler (1988) stated below.

THEOREM EC.5 (Winkler (1988) Theorem 2.1(b)). *Let \mathcal{X} be a measurable space with σ -field \mathcal{F} and suppose that \mathcal{P} is a simplex of probability measures whose extreme points are Dirac measures. Fix measurable functions f_1, \dots, f_n and real numbers c_1, \dots, c_n . Consider the set*

$$\mathcal{H} = \left\{ q \in \mathcal{P} : f_i \text{ is } q\text{-integrable and } \int f_i dq = c_i, \ 1 \leq i \leq n \right\}$$

Then \mathcal{H} is convex and

$$\text{ex } \mathcal{H} = \left\{ q \in \mathcal{H} : q = \sum_{i=1}^m t_i \cdot \delta(x_i), \ t_i > 0, \ \sum_{i=1}^m t_i = 1, \ x_i \in \mathcal{X}, \ 1 \leq m \leq n+1, \right. \\ \left. \text{the vectors } (f_1(x_i), \dots, f_n(x_i), 1), \ 1 \leq i \leq m, \text{ are linearly independent} \right\}$$

where $\text{ex } \mathcal{H}$ denotes the set of all extreme points of \mathcal{H} .

THEOREM EC.6 (Adapted from Winkler (1988) Theorem 3.2). *Let \mathcal{X} be a Hausdorff space, \mathcal{F} be the Borel σ -field and $\mathcal{P}_r(\mathcal{X})$ be the set of regular probability measures on \mathcal{X} . Let*

$$\mathcal{H} = \left\{ q \in \mathcal{P}_r(\mathcal{X}) : f_i \text{ is } q\text{-integrable and } \int f_i dq = c_i, \ 1 \leq i \leq n \right\}$$

Every measure affine functional J on \mathcal{H} fulfills

$$\sup\{J(q) : q \in \mathcal{H}\} = \sup\{J(q) : q \in \text{ex } \mathcal{H}\}$$

Theorem EC.6 is precisely Theorem 3.2 in Winkler (1988), except replacing the inequalities with equalities for the moments that define \mathcal{H} , which is immediate (and is pointed out by the comment right after the theorem in Winkler (1988)).

PROPOSITION EC.1 (Winkler (1988) Proposition 3.1). *Let \mathcal{X} , \mathcal{F} and \mathcal{H} be given as in Theorem EC.6 and the function g on \mathcal{X} be integrable for every $q \in \mathcal{H}$ (possibly with integral values ∞ or $-\infty$). Then the functional G on \mathcal{H} defined by $G(q) = \int_{\mathcal{X}} g dq$ is measure affine.*

Proof of Theorem 3. By Examples 2.1(a) in Winkler (1988), the set \mathcal{P} in Theorem EC.5 can be chosen to be the set of all regular probability measures. On Polish space every probability measure is regular. Therefore, on the space \mathbb{R}^+ , which is Polish, we can take \mathcal{P} in Theorem EC.5 as the set of all probability measures. The \mathcal{H} in Theorems EC.5 and EC.6 then coincide. By Proposition EC.1, the objective $\nu \mathbb{E}[H(X)]$ in $OPT(\mathcal{P}^+)$ is measure affine. Therefore, using Theorems EC.5 and EC.6, and noting that

$$\mathcal{H} \supseteq \left\{ q \in \mathcal{H} : q = \sum_{i=1}^m t_i \cdot \delta(x_i), \ t_i > 0, \ \sum_{i=1}^m t_i = 1, \ x_i \in \mathcal{X}, \ 1 \leq m \leq n+1 \right\} \supseteq_{\text{ex}} \mathcal{H}$$

for the coincided \mathcal{H} in Theorems EC.5 and EC.6, we conclude the theorem. \square

Proof of Proposition 1. If program (5) is consistent, then by Theorem 3, either an optimal solution in \mathcal{P}_3^+ exists, which corresponds to the first case of the lemma, or there exists a feasible sequence $\mathbb{P}^{(k)} \in \mathcal{P}_3^+$ such that $Z(\mathbb{P}^{(k)}) \rightarrow Z^*$. Let $\mathbb{P}^{(k)} \sim (x_1^{(k)}, x_2^{(k)}, x_3^{(k)}, p_1^{(k)}, p_2^{(k)}, p_3^{(k)})$. Suppose that x_i 's are all bounded above by a number, say M . Then, since $[0, M]^3 \times \mathcal{S}_3$ is a compact set, by Bolzano-Weierstrass Theorem we must have a subsequence of $(x_1^{(k)}, x_2^{(k)}, x_3^{(k)}, p_1^{(k)}, p_2^{(k)}, p_3^{(k)})$, say $(x_1^{(k_j)}, x_2^{(k_j)}, x_3^{(k_j)}, p_1^{(k_j)}, p_2^{(k_j)}, p_3^{(k_j)})$ converge to $(x_1^*, x_2^*, x_3^*, p_1^*, p_2^*, p_3^*)$ in $[0, M]^3 \times \mathcal{S}_3$. Since $H(x)$ is continuous by construction, we have $Z(\mathbb{P}^{(k_j)}) = \nu \sum_{i=1}^3 H(x_i^{(k_j)}) p_i^{(k_j)} \rightarrow \nu \sum_{i=1}^3 H(x_i^*) p_i^* = Z(\mathbb{P}^*)$, where $\mathbb{P}^* \sim (x_1^*, x_2^*, x_3^*, p_1^*, p_2^*, p_3^*)$. As $Z(\mathbb{P}^{(k_j)})$ is a subsequence of $Z(\mathbb{P}^{(k)})$, $Z(\mathbb{P}^*)$ must be equal to Z^* , and so \mathbb{P}^* is an optimal solution, which reduces to the first case in the lemma. Therefore, for the second case, we should focus on the scenario that at least one $x_i^{(k)}$ satisfies $\limsup_{k \rightarrow \infty} x_i^{(k)} = \infty$.

Without loss of generality, we fix the convention that $x_1^{(k)} \leq x_2^{(k)} \leq x_3^{(k)}$. If at least one of $x_i^{(k)}$ satisfies $\limsup_{k \rightarrow \infty} x_i^{(k)} = \infty$, we must have $\limsup_{k \rightarrow \infty} x_3^{(k)} = \infty$. In order that $\mathbb{P}^{(k)}$ is feasible, $\mathbb{E}^{(k)}[X] = \mu$ holds and so $x_1^{(k)} \leq \mu$ for all k . We now distinguish two cases: either $x_2^{(k)}$ is uniformly bounded, say by a large number $M \geq \mu$, or $\limsup_{k \rightarrow \infty} x_2^{(k)} = \infty$ also. Consider the first case. First,

we find a subsequence $x_3^{(k_j)} \nearrow \infty$. Since $(x_1^{(k_j)}, x_2^{(k_j)}) \in [0, M]^2$ which is compact, we can choose a further subsequence $k_{j'}$ such that $(x_1^{(k_{j'})}, x_2^{(k_{j'})}, x_3^{(k_{j'})}) \rightarrow (x_1^*, x_2^*, \infty)$ where $(x_1^*, x_2^*) \in [0, M]^2$. Now, since $(p_1^{(k_{j'})}, p_2^{(k_{j'})}, p_3^{(k_{j'})}) \in \mathcal{S}_3$ which is also compact, we can choose another further subsequence $k_{j''}$ such that $(p_1^{(k_{j''})}, p_2^{(k_{j''})}, p_3^{(k_{j''})}) \rightarrow (p_1^*, p_2^*, p_3^*) \in \mathcal{S}_3$. Note that by the constraint $\mathbb{E}^{(k)}[X^2] = p_1^{(k_{j''})} x_1^{(k_{j''})^2} + p_2^{(k_{j''})} x_2^{(k_{j''})^2} + p_3^{(k_{j''})} x_3^{(k_{j''})^2} = \sigma$, we must have $p_3^{(k_{j''})} = (\sigma - p_1^{(k_{j''})} x_1^{(k_{j''})^2} - p_2^{(k_{j''})} x_2^{(k_{j''})^2}) / x_3^{(k_{j''})^2} \leq \sigma / x_3^{(k_{j''})^2} \rightarrow 0$. In conclusion, in this case, we end up being able to find a sequence of measures $\mathbb{P}^{(k)'} \sim (x_1^{(k)'}, x_2^{(k)'}, x_3^{(k)'}, p_1^{(k)'}, p_2^{(k)'}, p_3^{(k)'})$ with $(x_1^{(k)'}, x_2^{(k)'}, x_3^{(k)'}, p_1^{(k)'}, p_2^{(k)'}, p_3^{(k)'}) \rightarrow (x_1^*, x_2^*, \infty, p_1^*, p_2^*, 0)$ where $x_1^*, x_2^* \in \mathbb{R}^+$ and $(p_1^*, p_2^*) \in \mathcal{S}_2$.

For the second case, namely when $\limsup_{k \rightarrow \infty} x_i^{(k)} = \infty$ for both $i = 2$ and 3 . We can argue similarly that there is a sequence of measures $\mathbb{P}^{(k)'} \sim (x_1^{(k)'}, x_2^{(k)'}, x_3^{(k)'}, p_1^{(k)'}, p_2^{(k)'}, p_3^{(k)'})$, such that $x_2^{(k)'}, x_3^{(k)'} \rightarrow \infty$ and $p_2^{(k)'}, p_3^{(k)'} \rightarrow 0$. In other words, $(x_1^{(k)'}, x_2^{(k)'}, x_3^{(k)'}, p_1^{(k)'}, p_2^{(k)'}, p_3^{(k)'}) \rightarrow (x_1^*, \infty, \infty, 1, 0, 0)$ where $x_1^* \in \mathbb{R}^+$. \square

Proof of Lemma 2. It follows from Jensen's inequality that for any $\mathbb{P} \in \mathcal{P}^+$, $\mathbb{E}[X^2] \geq \mathbb{E}[X]^2$, which gives $\sigma \geq \mu^2$ in (5). On the other hand, if $\sigma \geq \mu^2$, it is also rudimentary to find $\mathbb{P} \in \mathcal{P}_2^+$ with $\mathbb{E}[X] = \mu$ and $\mathbb{E}[X^2] = \sigma$. Substituting $\mu = \eta/\nu$ and $\sigma = 2\beta/\nu$, we get $\eta^2 \leq 2\beta\nu$. Lastly, $\mathbb{E}[X^2] = \mathbb{E}[X]^2$ if and only if \mathbb{P} is a point mass. The equivalent statements regarding program (1) follows from Theorem 2. \square

Proof of Proposition 2. Consider a sequence $f^{(k)}(x), x \geq a$ given by

$$f^{(k)}(x) = \begin{cases} \eta - \nu(x - a) & \text{for } a \leq x \leq x_1^{(k)} + a \\ \eta - \nu x_1^{(k)} - \nu p_2^{(k)}(x - a - x_1^{(k)}) & \text{for } x_1^{(k)} + a \leq x \leq x_2^{(k)} + a \\ 0 & \text{for } x_2^{(k)} + a \leq x \end{cases}$$

where

$$\begin{aligned} x_1^{(k)} &= \mu - \gamma^{(k)} \text{ and } \gamma^{(k)} = \frac{\sigma - \mu^2}{x_2^{(k)} - \mu} \\ x_2^{(k)} &\rightarrow \infty \\ p_1^{(k)} &= 1 - p_2^{(k)} \\ p_2^{(k)} &= \frac{\sigma - \mu^2}{x_2^{(k)^2} - 2\mu x_2^{(k)} + \sigma} \end{aligned} \tag{EC.9}$$

and μ, σ are defined in (4).

It is obvious that for large enough $x_2^{(k)}$, $f^{(k)}$ is non-negative and convex. Moreover, $f^{(k)}(a) = f^{(k)}(a+) = \eta$ and $f_+^{(k)'}(a) \geq -\nu$. To show $\int_a^\infty f(x)dx = \beta$, we first verify that

$$p_1^{(k)} x_1^{(k)} + p_2^{(k)} x_2^{(k)} = \mu \quad (\text{EC.10})$$

and

$$p_1^{(k)} x_1^{(k)2} + p_2^{(k)} x_2^{(k)2} = \sigma \quad (\text{EC.11})$$

for all k . In fact, we will do so by showing that $\gamma^{(k)}$ and $p_2^{(k)}$ displayed in (EC.9) are the unique choices that satisfy (EC.10) and (EC.11) and also $x_1^{(k)} = \mu - \gamma^{(k)}$ and $p_1^{(k)} = 1 - p_2^{(k)}$. With the latter conditions, (EC.10) and (EC.11) can be written as

$$(1 - p_2^{(k)})(\mu - \gamma^{(k)}) + p_2^{(k)} x_2^{(k)} = \mu$$

and

$$(1 - p_2^{(k)})(\mu - \gamma^{(k)})^2 + p_2^{(k)} x_2^{(k)2} = \sigma$$

respectively, which further gives

$$p_2^{(k)} (\gamma^{(k)} + x_2^{(k)} - \mu) - \gamma^{(k)} = 0 \quad (\text{EC.12})$$

and

$$p_2^{(k)} \left(x_2^{(k)2} - (\mu - \gamma^{(k)})^2 \right) + (\mu - \gamma^{(k)})^2 = \sigma \quad (\text{EC.13})$$

From (EC.12) we have

$$p_2^{(k)} = \frac{\gamma^{(k)}}{\gamma^{(k)} + x_2^{(k)} - \mu} \quad (\text{EC.14})$$

Putting (EC.14) into (EC.13), we get

$$\frac{\gamma^{(k)}}{\gamma^{(k)} + x_2^{(k)} - \mu} \left(x_2^{(k)2} - (\mu - \gamma^{(k)})^2 \right) + (\mu - \gamma^{(k)})^2 = \sigma$$

which can be simplified to

$$\gamma^{(k)} \left(x_2^{(k)} + \mu - \gamma^{(k)} \right) + (\mu - \gamma^{(k)})^2 = \sigma$$

giving

$$\gamma^{(k)} = \frac{\sigma - \mu^2}{x_2^{(k)} - \mu} \quad (\text{EC.15})$$

Plugging (EC.15) into (EC.14), we have

$$p_2^{(k)} = \frac{\sigma - \mu^2}{(\sigma - \mu^2) + (x_2^{(k)} - \mu)^2} \quad (\text{EC.16})$$

thus recovering $\gamma^{(k)}$ and $p_2^{(k)}$ in (EC.9).

Therefore,

$$\begin{aligned} \int_a^\infty f^{(k)}(x)dx &= \int_a^{x_1^{(k)}+a} [\eta - \nu(x-a)]dx + \int_{x_1^{(k)}+a}^{x_2^{(k)}+a} [\eta - \nu x_1^{(k)} - \nu p_2^{(k)}(x-a-x_1^{(k)})]dx \\ &= \eta x_2^{(k)} - \frac{\nu}{2} x_1^{(k)2} - \nu x_1^{(k)}(x_2^{(k)} - x_1^{(k)}) - \frac{\nu p_2^{(k)}}{2} (x_2^{(k)} - x_1^{(k)})^2 \\ &= \eta x_2^{(k)} + \frac{\nu p_1^{(k)}}{2} x_1^{(k)2} + \frac{\nu p_2^{(k)}}{2} x_2^{(k)2} - \nu p_1^{(k)} x_1^{(k)} x_2^{(k)} - \nu p_2^{(k)} x_2^{(k)2} \text{ using } p_1^{(k)} = 1 - p_2^{(k)} \\ &= \eta x_2^{(k)} + \frac{\nu \sigma}{2} - \nu x_2^{(k)} \mu \text{ using (EC.10) and (EC.11)} \\ &= \beta \text{ using } \eta - \nu \mu = 0 \text{ and } \beta = \nu \sigma / 2 \end{aligned}$$

Hence $f^{(k)}$ is feasible for (1) for large enough k .

Now, the objective value evaluated at $f^{(k)}$ is

$$\int_a^{x_1^{(k)}+a} h(x)(\eta - \nu(x-a))dx + \int_{x_1^{(k)}+a}^{x_2^{(k)}+a} h(x)(\eta - \nu x_1^{(k)} - \nu p_2^{(k)}(x-a-x_1^{(k)}))dx \quad (\text{EC.17})$$

The first term in (EC.17) is bounded since $x_1^{(k)} \rightarrow \mu$. We focus on the second term. By the assumption, we can find $C > 0$ such that $h(x) \geq Cx^\epsilon$ for all $x \geq a$. Then, for large enough k ,

$$\begin{aligned} &\int_{x_1^{(k)}+a}^{x_2^{(k)}+a} h(x)(\eta - \nu x_1^{(k)} - \nu p_2^{(k)}(x-a-x_1^{(k)}))dx \\ &\geq C \int_{x_1^{(k)}+a}^{x_2^{(k)}+a} x^\epsilon (\eta - \nu x_1^{(k)} - \nu p_2^{(k)}(x-a-x_1^{(k)}))dx \\ &\geq C \int_{x_1^{(k)}+a}^{x_2^{(k)}+a} [(\eta - \nu p_1^{(k)} x_1^{(k)} + \nu p_2^{(k)} a) x^\epsilon - \nu p_2^{(k)} x^{\epsilon+1}]dx \\ &= (\eta - \nu p_1^{(k)} x_1^{(k)} + \nu p_2^{(k)} a) \frac{x^{\epsilon+1}}{\epsilon+1} \Big|_{x_1^{(k)}+a}^{x_2^{(k)}+a} - \nu p_2^{(k)} \frac{x^{\epsilon+2}}{\epsilon+2} \Big|_{x_1^{(k)}+a}^{x_2^{(k)}+a} \\ &= (\eta - \nu p_1^{(k)} x_1^{(k)} + \nu p_2^{(k)} a) \frac{(x_2^{(k)}+a)^{\epsilon+1}}{\epsilon+1} - (\eta - \nu p_1^{(k)} x_1^{(k)} + \nu p_2^{(k)} a) \frac{(x_1^{(k)}+a)^{\epsilon+1}}{\epsilon+1} \\ &\quad - \nu p_2^{(k)} \frac{(x_2^{(k)}+a)^{\epsilon+2}}{\epsilon+2} + \nu p_2^{(k)} \frac{(x_1^{(k)}+a)^{\epsilon+2}}{\epsilon+2} \end{aligned} \quad (\text{EC.18})$$

Note that since $p_1^{(k)} \rightarrow 1$, $x_1^{(k)} \rightarrow \mu$, $p_2^{(k)} \rightarrow 0$ and $\eta - \nu\mu = 0$, the second term in (EC.18) converges to 0. Moreover, since $p_2^{(k)} \rightarrow 0$, the fourth term also converges to 0. Consider the first term in (EC.18).

In particular,

$$\begin{aligned} \eta - \nu p_1^{(k)} x_1^{(k)} + \nu p_2^{(k)} a &= \eta - \nu(1 - p_2^{(k)})(\mu - \gamma^{(k)}) + \nu p_2^{(k)} a \\ &= p_1^{(k)} \nu \gamma^{(k)} + \nu p_2^{(k)} (\mu + a) \end{aligned}$$

by using $\eta - \nu\mu = 0$ and $p_1^{(k)} = 1 - p_2^{(k)}$. Substituting $\gamma^{(k)} = (\sigma - \mu^2)/(x_2^{(k)} - \mu)$ and $p_2^{(k)} = \Theta(1/x_2^{(k)2})$, and using $p_1^{(k)} \rightarrow 1$, we have

$$(\eta - \nu p_1^{(k)} x_1^{(k)} + \nu p_2^{(k)} a) \frac{(x_2^{(k)} + a)^{\epsilon+1}}{\epsilon + 1} = (p_1^{(k)} \nu \gamma^{(k)} + \nu p_2^{(k)} (\mu + a)) \frac{(x_2^{(k)} + a)^{\epsilon+1}}{\epsilon + 1} = \frac{\nu(\sigma - \mu^2) x_2^{(k)\epsilon}}{\epsilon + 1} (1 + o(1))$$

On the other hand, for the third term in (EC.18), substituting $p_2^{(k)} = (\sigma - \mu^2)/(x_2^{(k)2} - 2\mu x_2^{(k)} + \sigma)$, we have

$$-\nu p_2^{(k)} \frac{(x_2^{(k)} + a)^{\epsilon+2}}{\epsilon + 2} = -\frac{\nu(\sigma - \mu^2) x_2^{(k)\epsilon}}{\epsilon + 2} (1 + o(1))$$

Thus, (EC.18) is equal to

$$\left(\frac{1}{\epsilon + 1} - \frac{1}{\epsilon + 2} \right) \nu(\sigma - \mu^2) x_2^{(k)\epsilon} (1 + o(1)) \rightarrow \infty$$

and hence the optimal value of (1) is ∞ . □

EC.2. Proofs for Section 5

To prove Proposition 3, we borrow the following result:

LEMMA EC.2 (Adapted from Theorem 5.1 in Birge and Dulá (1991)). *Consider*

$OPT(\mathcal{P}[0, \tilde{c}])$ for any $0 < \tilde{c} < \infty$. Suppose H is convex with derivative H' convex on $(0, c)$ and concave on (c, \tilde{c}) for some $0 \leq c \leq \tilde{c}$. If $OPT(\mathcal{P}[0, \tilde{c}])$ is consistent, then an optimal solution exists and lies in $\mathcal{P}_2[0, \tilde{c}]$.

This lemma follows from Theorem 5.1 in Birge and Dulá (1991) that applies to the associated dual problem.

Proof of Proposition 3. By Theorem 3, $OPT(\mathcal{P}^+)$ has the same optimal value as $OPT(\mathcal{P}_3^+)$. By Lemma EC.2, for every $\mathbb{P} \in \mathcal{P}_3^+$, which necessarily has bounded support say on $[0, M]$ for some $M > 0$, there exists $\mathbb{P}' \in \mathcal{P}_2[0, M]$ such that $Z(\mathbb{P}') \geq Z(\mathbb{P})$. Hence $OPT(\mathcal{P}_3^+)$ has the same optimal value as $OPT(\mathcal{P}_2^+)$, which concludes the proposition. \square

Proof of Proposition 4. Proof of 1: Let the optimal probability measure in \mathcal{P}_2^+ be represented by (x_1, x_2, p_1, p_2) . Note that $x_1 \neq x_2$ since otherwise $\sigma = \mu^2$. Adopting a similar line of analysis as in Birge and Dulá (1991), we let $x_1 < x_2$ without loss of generality. For a two-support-point distribution to be feasible, we must have $x_1 < \mu$. Feasibility also enforces that $p_1 x_1 + p_2 x_2 = \mu$, $p_1 x_1^2 + p_2 x_2^2 = \sigma$ and $p_1 + p_2 = 1$. Hence $p_2 = 1 - p_1$, which gives $p_1 x_1 + (1 - p_1) x_2 = \mu$ and $p_1 x_1^2 + (1 - p_1) x_2^2 = \sigma$. From the first equation we get $p_1 = (x_2 - \mu)/(x_2 - x_1)$. Putting this into $p_1 x_1^2 + (1 - p_1) x_2^2 = \sigma$, we further get $x_2 = (\sigma - \mu x_1)/(\mu - x_1)$. Now, putting this in turn into $p_1 = (x_2 - \mu)/(x_2 - x_1)$, we obtain $p_1 = (\sigma - \mu^2)/(\sigma - 2\mu x_1 + x_1^2)$ and hence $p_2 = 1 - p_1 = (\mu - x_1)^2/(\sigma - 2\mu x_1 + x_1^2)$. Therefore, Z^* is given by

$$\max_{x_1 \in [0, \mu)} \nu(p_1 H(x_1) + p_2 H(x_2)) = \max_{x_1 \in [0, \mu)} \nu \left(\frac{\sigma - \mu^2}{\sigma - 2\mu x_1 + x_1^2} H(x_1) + \frac{(\mu - x_1)^2}{\sigma - 2\mu x_1 + x_1^2} H\left(\frac{\sigma - \mu x_1}{\mu - x_1}\right) \right)$$

which is exactly $\max_{x_1 \in [0, \mu)} W(x_1)$.

Proof of 2: Let $\mathbb{P}^{(k)} \sim (x_1^{(k)}, x_2^{(k)}, p_1^{(k)}, p_2^{(k)})$ be a feasible sequence with $Z(\mathbb{P}^{(k)}) \rightarrow Z^*$. Without loss of generality let $x_1^{(k)} \leq x_2^{(k)}$. Since $p_1^{(k)} x_1^{(k)} + p_2^{(k)} x_2^{(k)} = \mu$, we must have $x_1^{(k)} \leq \mu$. Then we must have a subsequence $x_2^{(k_i)} \rightarrow \infty$, since otherwise $(x_1^{(k)}, x_2^{(k)}, p_1^{(k)}, p_2^{(k)})$ would lie in a compact set and there would exist a subsequence $(x_1^{(k'_i)}, x_2^{(k'_i)}, p_1^{(k'_i)}, p_2^{(k'_i)}) \rightarrow (x_1^*, x_2^*, p_1^*, p_2^*)$, where $Z(\mathbb{P}^{(k'_i)}) = \nu \sum_{j=1}^2 p_j^{(k'_i)} H(x_j^{(k'_i)}) \rightarrow \nu \sum_{j=1}^2 p_j^* H(x_j^*)$ by the continuity of H , violating the non-existence of optimal solution. By $p_1^{(k_i)} x_1^{(k_i)^2} + p_2^{(k_i)} x_2^{(k_i)^2} = \sigma$, we have $p_2^{(k_i)} = (\sigma - p_1^{(k_i)} x_1^{(k_i)^2})/x_2^{(k_i)^2} \rightarrow 0$, and $p_2^{(k_i)} x_2^{(k_i)} = (\sigma - p_1^{(k_i)} x_1^{(k_i)^2})/x_2^{(k_i)} \rightarrow 0$. Thus $p_1^{(k_i)} = 1 - p_2^{(k_i)} \rightarrow 1$ and $x_1^{(k_i)} = (\mu - p_2^{(k_i)} x_2^{(k_i)})/p_1^{(k_i)} \rightarrow \mu$. Therefore,

$$\begin{aligned} Z(\mathbb{P}^{(k_i)}) &= \nu \left(p_1^{(k_i)} H(x_1^{(k_i)}) + p_2^{(k_i)} H(x_2^{(k_i)}) \right) = \nu \left(p_1^{(k_i)} H(x_1^{(k_i)}) + \frac{\sigma - p_1^{(k_i)} x_1^{(k_i)^2}}{x_2^{(k_i)^2}} H(x_2^{(k_i)}) \right) \\ &\rightarrow \nu(H(\mu) + \lambda(\sigma - \mu^2)) \end{aligned}$$

Proof of 3: First, we show that $W(x_1) \rightarrow \nu(H(\mu) + \lambda(\sigma - \mu^2))$ as $x_1 \nearrow \mu$. Consider the second term of $W(x_1)$ given by

$$\lim_{x_1 \nearrow \mu} \frac{\nu(\mu - x_1)^2}{\sigma - 2\mu x_1 + x_1^2} H\left(\frac{\sigma - \mu x_1}{\mu - x_1}\right) = \lim_{x_1 \nearrow \mu} \frac{\nu(\sigma - \mu x_1)^2}{\sigma - 2\mu x_1 + x_1^2} \left(\frac{\mu - x_1}{\sigma - \mu x_1}\right)^2 H\left(\frac{\sigma - \mu x_1}{\mu - x_1}\right) = \nu\lambda(\sigma - \mu^2)$$

and the claim follows. Combining Parts 1 and 2 of this proposition, we must have $Z^* = \max_{x_1 \in [0, \mu]} W(x_1)$. \square

EC.3. Proofs for Section 6

We first show a result in parallel to Theorem 2 for the case of (11):

THEOREM EC.7. *Suppose h is bounded. Then the optimal value of (11) is the same as*

$$\begin{aligned} \max_{\mathbb{P}} \quad & \bar{\nu} \mathbb{E}[H(X)] \\ \text{subject to} \quad & \underline{\mu} \leq \mathbb{E}[X] \leq \bar{\mu} \\ & \underline{\sigma} \leq \mathbb{E}[X^2] \leq \bar{\sigma} \\ & \mathbb{P} \in \mathcal{P}^+ \end{aligned} \tag{EC.19}$$

Here the decision variable is a probability distribution $\mathbb{P} \in \mathcal{P}^+$, and $\mathbb{E}[\cdot]$ is the corresponding expectation. Moreover, there is a one-to-one correspondence between the feasible solutions to (11) and (EC.19), given by $f'_+(x+a) = \bar{\nu}(p(x) - 1)$ for $x \in \mathbb{R}^+$, where f'_+ is the right derivative of a feasible solution f of (11) such that $f(x) = \int_a^x f'_+(t)dt + \eta$ for $x \geq a$, and p is a probability distribution function that is associated with a feasible probability measure over \mathbb{R}^+ in (EC.19).

Proof of Theorem EC.7. Note that formulation (11) can be written as

$$\begin{aligned} \max_{\underline{\beta} \leq \beta \leq \bar{\beta}, \underline{\eta} \leq \eta \leq \bar{\eta}} \max_f \quad & \int_a^\infty h(x) f(x) dx \\ \text{subject to} \quad & \int_a^\infty f(x) dx = \beta \\ & f(a) = f(a+) = \eta \\ & f'_+(a) \geq -\bar{\nu} \\ & f(x) \text{ convex for } x \geq a \\ & f(x) \geq 0 \text{ for } x \geq a \end{aligned} \tag{EC.20}$$

The inner maximization is exactly (1), and thus by Theorem 2 we can reformulate (EC.20) as

$$\begin{aligned} \max_{\underline{\beta} \leq \beta \leq \bar{\beta}, \underline{\eta} \leq \eta \leq \bar{\eta}} \max_{\mathbb{P}} \quad & \bar{\nu} \mathbb{E}[H(X)] \\ \text{subject to} \quad & \mathbb{E}[X] = \frac{\eta}{\bar{\nu}} \\ & \mathbb{E}[X^2] = \frac{2\beta}{\bar{\nu}} \\ & \mathbb{P} \in \mathcal{P}^+ \end{aligned}$$

which is equivalent to (EC.19). □

For convenience, we denote $\widetilde{OPT}(\mathcal{D})$ as the program

$$\begin{aligned} \max_{\mathbb{P}} \quad & \bar{\nu} \mathbb{E}[H(X)] \\ \text{subject to} \quad & \underline{\mu} \leq \mathbb{E}[X] \leq \bar{\mu} \\ & \underline{\sigma} \leq \mathbb{E}[X^2] \leq \bar{\sigma} \\ & \mathbb{P} \in \mathcal{D} \end{aligned}$$

where \mathcal{D} is a collection of probability measures on \mathbb{R} . For example, (EC.19) can be written as $\widetilde{OPT}(\mathcal{P}^+)$. Let $\tilde{Z}(\mathbb{P}) = \bar{\nu} \mathbb{E}[H(X)]$ be the objective function in \mathbb{P} .

PROPOSITION EC.2. *The optimal value of $\widetilde{OPT}(\mathcal{P}^+)$ is identical to that of $\widetilde{OPT}(\mathcal{P}_3^+)$.*

Proof of Proposition EC.2. For \mathbb{P} feasible in $\widetilde{OPT}(\mathcal{P}^+)$, let $\mu = \mathbb{E}[X]$ and $\sigma = \mathbb{E}[X^2]$ be its first and second moments. By Theorem 3 there must exist $\mathbb{P}' \in \mathcal{P}_3^+$ with the corresponding expectations $\mathbb{E}'[X] = \mu$ and $\mathbb{E}'[X^2] = \sigma$ such that $\tilde{Z}(\mathbb{P}) \leq \tilde{Z}(\mathbb{P}')$. □

Proof of Theorem 5. Theorem 5 follows from Theorem EC.7 and Proposition EC.2, in the same way as the proof of Theorem 1. □

PROPOSITION EC.3. *Under Assumption 3, $\widetilde{OPT}(\mathcal{P}^+)$ has the same optimal value as $\widetilde{OPT}(\mathcal{P}_2^+)$.*

Proof of Proposition EC.3. We know from Proposition EC.2 that $\widetilde{OPT}(\mathcal{P}^+)$ has the same optimal value as $\widetilde{OPT}(\mathcal{P}_3^+)$. Any $\mathbb{P} \in \mathcal{P}_3^+$ must necessarily have bounded support, say on $[0, M]$. By Lemma EC.2 there must exist $\mathbb{P}' \in \mathcal{P}_2^+$, with the same first and second moments as \mathbb{P} , such that $\tilde{Z}(\mathbb{P}) \leq \tilde{Z}(\mathbb{P}')$. □

The following explains the origin of the two subproblems in (13):

LEMMA EC.3. *Under Assumption 1, and let $\bar{\sigma} \geq \underline{\mu}^2$. The optimal value of $\widetilde{OPT}(\mathcal{P}_2^+)$ is given by $\tilde{Z}^* = \max\{\tilde{Z}_1^*, \tilde{Z}_2^*\}$, where \tilde{Z}_1^* is the optimal value of*

$$\begin{aligned}
 & \max_{\mathbb{P}} \quad \bar{\nu} \mathbb{E}[H(X)] \\
 & \text{subject to} \quad \mathbb{E}[X] = \bar{\mu} \\
 & \quad \underline{\sigma} \leq \mathbb{E}[X^2] \leq \bar{\sigma} \\
 & \quad \mathbb{P} \in \mathcal{P}_2^+
 \end{aligned} \tag{EC.21}$$

and \tilde{Z}_2^* is the optimal value of

$$\begin{aligned}
 & \max_{\mathbb{P}} \quad \bar{\nu} \mathbb{E}[H(X)] \\
 & \text{subject to} \quad \underline{\mu} \leq \mathbb{E}[X] \leq \bar{\mu} \\
 & \quad \mathbb{E}[X^2] = \bar{\sigma} \\
 & \quad \mathbb{P} \in \mathcal{P}_2^+
 \end{aligned} \tag{EC.22}$$

respectively.

Proof of Lemma EC.3. We argue that to solve $\widetilde{OPT}(\mathcal{P}_2^+)$, it suffices to restrict attention to the feasible region $\{\mathbb{P} \in \mathcal{P}_2^+ : \mathbb{E}[X] = \bar{\mu}, \underline{\sigma} \leq \mathbb{E}[X^2] \leq \bar{\sigma}\} \cup \{\mathbb{P} \in \mathcal{P}_2^+ : \underline{\mu} \leq \mathbb{E}[X] \leq \bar{\mu}, \mathbb{E}[X^2] = \bar{\sigma}\}$. Since $h \geq 0$, $\tilde{Z}^* \geq 0$. There is nothing to prove if $\tilde{Z}^* = 0$. So suppose $\tilde{Z}^* > 0$. There exists $\mathbb{P} \sim (x_1, x_2, p_1, p_2) \in \mathcal{P}_2^+$ with one of the x_i 's having $H(x_i) > 0$ and $p_i > 0$. Now suppose \mathbb{P} satisfies $\mathbb{E}[X] < \bar{\mu}$ and $\mathbb{E}[X^2] < \bar{\sigma}$. We can increase x_i so that $\mathbb{E}[X] \leq \bar{\mu}$ and $\mathbb{E}[X^2] \leq \bar{\sigma}$ remain satisfied, and $\tilde{Z}^*(\mathbb{P})$ is at least as large as before since $H(x)$ is non-decreasing. Hence any \mathbb{P} such that $\mathbb{E}[X] < \bar{\mu}$ and $\mathbb{E}[X^2] < \bar{\sigma}$ must have $\tilde{Z}(\mathbb{P}) \leq \tilde{Z}(\mathbb{P}')$ for some $\mathbb{P}' \in \{\mathbb{P} \in \mathcal{P}_2^+ : \mathbb{E}[X] = \bar{\mu}, \underline{\sigma} \leq \mathbb{E}[X^2] \leq \bar{\sigma}\} \cup \{\mathbb{P} \in \mathcal{P}_2^+ : \underline{\mu} \leq \mathbb{E}[X] \leq \bar{\mu}, \mathbb{E}[X^2] = \bar{\sigma}\}$. This proves the lemma. \square

Proof of Theorem 6. Lemma EC.3 allows one to consider only the programs (EC.21) and (EC.22) when solving $\widetilde{OPT}(\mathcal{P}_2^+)$. Theorem 6 then follows from Lemma 3, Theorem EC.7 and Proposition EC.3, using the same line of arguments in the proof of Theorem 4. \square